

Fast Graph Generation via Spectral Diffusion

Tianze Luo[†], Zhanfeng Mo[†], Sinno Jialin Pan^{*}

Abstract—Generating graph-structured data is a challenging problem, which requires learning the underlying distribution of graphs. Various models such as graph VAE, graph GANs, and graph diffusion models have been proposed to generate meaningful and reliable graphs, among which the diffusion models have achieved state-of-the-art performance. In this paper, we argue that running full-rank diffusion SDEs on the whole graph adjacency matrix space hinders diffusion models from learning graph topology generation, and hence significantly deteriorates the quality of generated graph data. To address this limitation, we propose an efficient yet effective Graph Spectral Diffusion Model (GSDM), which is driven by low-rank diffusion SDEs on the graph spectrum space. Our spectral diffusion model is further proven to enjoy a substantially stronger theoretical guarantee than standard diffusion models. Extensive experiments across various datasets demonstrate that, our proposed GSDM turns out to be the SOTA model, by exhibiting both significantly higher generation quality and much less computational consumption than the baselines.

Index Terms—Graph generative model, graph diffusion, stochastic differential equations.

挑战

1 INTRODUCTION

LEARNING to generate graph-structural data not only requires knowing the nodes' feature distribution, but also a deep understanding of the underlying graph topology, which is essential to modelling various graph instances, such as social networks [1], [2], molecule structures [3], [4], neural architectures [5], recommender systems [6], etc. Conventional likelihood-based graph generative models, e.g. GraphGAN [7], GraphVAE [8] and GraphRNN [9], have demonstrated great strength on graph generation tasks. In general, a likelihood-based model is designed to learn the likelihood function of the underlying graph data distribution, with which one can draw new samples with preserved graph properties from the distribution of interest. However, most likelihood-based generative models suffer from either limited quality of modeling graph structures, or considerable computational burden [10].

Recently, a series of diffusion-based generative models have been proposed to overcome the limitations of likelihood-based models. Although being originally established for image generation [11], diffusion models exhibit a great success in graph generation tasks with complex graph structural properties [10], [12]. Roughly speaking, diffusion refers to a Stochastic Differential Equation (SDE) that smoothly converts authentic data into pure noise via noise insertion. A diffusion model is able to reverse the data back from the noise, as long as the corresponding reversed-time SDE is known, which depends on the time-dependent gradient field of the density function (i.e. score function). To achieve this, a neural network is trained to learn the score function, and hence a reversed-time process can be constructed with the estimated score. The first graph diffusion model through SDEs, coined Graph Diffusion via SDE Systems (GDSS) [10], is designed to simultaneously

generate node features and adjacency matrix via reversed diffusion. Similar to image diffusion models [11], [13], at each diffusion step, GDSS directly inserts standard Gaussian noise to both node features and the adjacency matrix. Meanwhile, two separate neural networks are trained to learn the score functions of the node features and adjacency matrix, respectively.

However, unlike the densely distributed image data, graph adjacency matrices can be highly sparse, which makes isotropic Gaussian noise insertion incompatible with graph structural data. In these circumstances, as shown in Figure 1, there is a stark difference between the diffusion process on images and on graph adjacency matrices. As can be seen from the figure, the image corrupted by full-rank Gaussian noise exhibits recognizable numerical patterns along the early- and middle-stage of forward diffusion. However, the corrupted sparse graph adjacency matrix degenerates into a dense matrix with uniformly distributed entries in a few diffusion steps. In intuition, Figure 1 implies that standard diffusion SDEs with full-rank isotropic noise insertion is destructive of learning graph topology and feature representations. Theoretically speaking, for extremely sparse graphs (e.g. molecules) with low-rank adjacency matrices, the adjacency score functions are supported on a low-dimensional manifold embedded in the full adjacency matrix space. Thus, directly applying diffusion models on graph topology generation is not desirable: once the diffusion SDE is run in the full space of the adjacency matrix, lethal noise will be injected into the out-of-support regions and drives the signal-to-noise ratio to be essentially zero, which is fatal for training score networks.

Even for densely connected graphs, the standard diffusion model is problematic for topology generation. Unlike image pixels that are merely locally correlated, an adjacency matrix governs the message-passing pattern of the whole graph. Thus, isotropic Gaussian noise insertion severely distorts the message-passing pattern, by blindly encouraging message passing on sparsely connected parts, which impedes the representation learning of sparse regions.

• T. Luo, Z. Mo and S. J. Pan are with the School of Computer Science and Engineering, Nanyang Technological University, Singapore. [†] indicates co-first authors with equal contribution. ^{*} indicates corresponding author. E-mail: {tianze001,zhanfeng001,sinnopan}@ntu.edu.sg

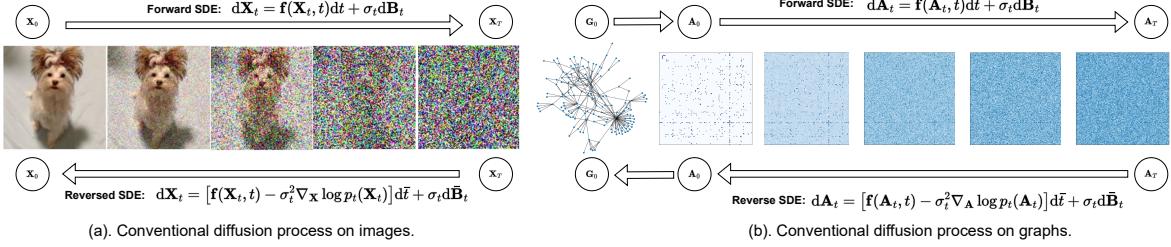


Fig. 1: Illustration of the difference between applying the conventional SDE diffusion process on images (shown in (a)) and on graphs (shown in (b)).

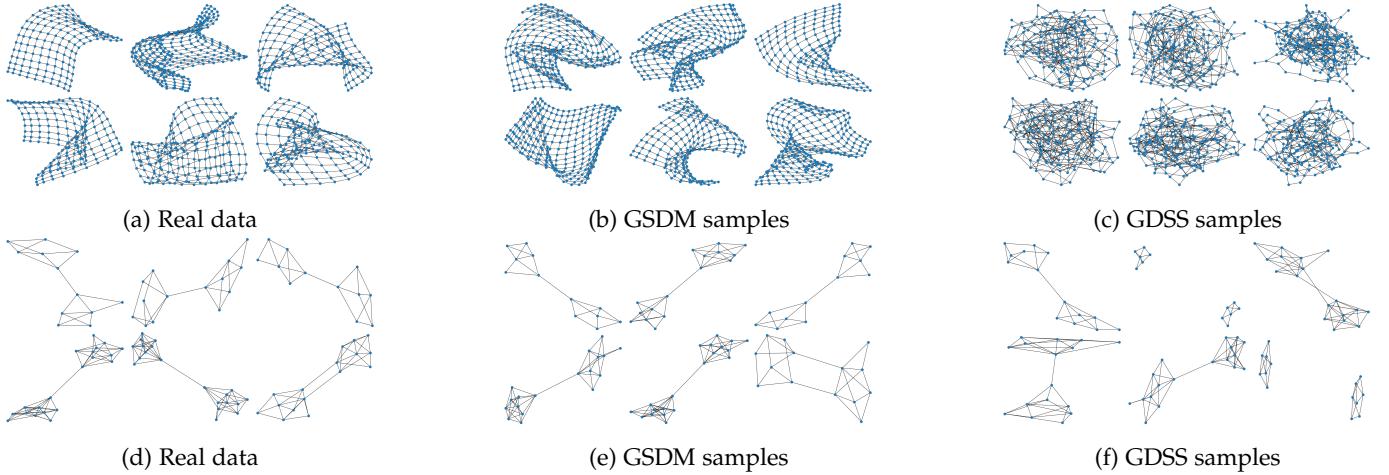


Fig. 2: Non-cherry-picked random samples from the testing set as well as samples generated by GSDM (ours) and GDSS [10], on Grid (top row) and Community-small (bottom row) datasets. For GDSS, we use the authors’ released code and checkpoints to generate the samples.

In order to establish a graph-friendly diffusion model, one should design an appropriate diffusion scheme that is compatible with the graph topology structure. To this end, we propose the Graph Spectral Diffusion Model (GSDM), which is driven by diffusion SDEs on both the node feature space and the graph spectrum space. At each diffusion step, instead of corrupting the entire adjacency matrix, our method confines the Gaussian insertion to the graph spectrum space, i.e. the eigenvalue matrix of the adjacency matrix. This novel diffusion scheme enables us to perform smooth transformations on graph data during both the training and sampling phases. As illustrated in Figure 2, our proposed GSDM significantly outperforms the standard graph diffusion model (GDSS [10]) in terms of graph generation quality and plausibility. For the Grid dataset shown in the top row of the figure, GDSS samples seem to be merely chaotic clusters, while GSDM samples exhibit smooth surface-like patterns that are visually similar to real data. For the Community-small dataset shown in the bottom row of the figure, GDSS fails to capture the link between two communities on some samples, while GSDM is able to capture dumbbell-like pattern (two clusters connected by one edge) as well as butterfly-like pattern (two clusters connected by two edges). This implies that GDSS’s generation not only fails to mimic the observed topology distribution, but also suffers from capturing challenging

details of the data such as links between communities. In contrast, GSDM is capable of generating high-quality graphs that are topologically similar to the real data, while retaining critical details.

We empirically evaluate the capability of our proposed GSDM on generic graph generation tasks, by evaluating the generation quality on both synthetic and real-world graph datasets. As shown in Section 4, GSDM outperforms existing one-shot generative models on various datasets, while achieving competitive performance to autoregressive models. Further molecule generation experiments show that our GSDM outperforms the state-of-the-art baselines, demonstrating that our proposed spectral diffusion model is capable of capturing complicated dependency between nodes and edges. Our main contributions are 3 folds:

- We propose a novel Graph Spectral Diffusion Model (GSDM) for fast and quality graph generation. Our method overcomes the limitations of existing graph diffusion models by leveraging diffusion SDEs on both the node feature and graph spectrum spaces. 创新点
- Through the lens of stochastic analysis, we prove that GSDM enjoys a substantially stronger performance guarantee than the standard graph diffusion model. Our proposed spectral diffusion sharpens the reconstruction error bound from $\mathcal{O}(n^2 \exp(n^2))$ to $\mathcal{O}(n \exp(n))$, where n is the number of nodes.

- We evaluate GSDM on both synthetic and real-world graph generation tasks, on which GSDM outperforms all existing graph generative models. Moreover, GSDM also achieves evidently higher computational efficiency compared to existing graph diffusion models.

2 RELATED WORK

From the view of graph generation strategies, graph generative models can be classified into two categories, sequential models and one-shot models. The sequential models, including GraphRNN [9], GraphVAE [8], generate the nodes and edges in a sequential way, i.e. one after another [14], with validity checks among the generation steps. Models using a similar generation process are also known as autoregressive models. In contrast, one-shot generative models, e.g. GAN-based models [15], VAE-based models [16], flow-based models [3] and score-based models [10], [12], learn the intrinsic dependency of graph structure by treating the distribution of all the components of a graph as a whole and therefore generate the entire graph in an integrative way. As a result, one-shot generative models enjoy the node permutation-invariant property. Moreover, one-shot models usually exhibit higher computational efficiency than autoregressive models. Our proposed GSDM generates data by reversing a spectral diffusion SDE with a learned spectral score function, thus it is essentially a score-based model. Since GSDM accepts the destination of the reversed diffusion SDE as the ultimately generated samples without any additional refinements, it can also be categorized as a one-shot model.

A recently proposed score-based model, GDSS [10], is the first and state-of-the-art diffusion-based generative model that simultaneously conducts nodes and edges generation. In essence, GDSS recasts the image diffusion paradigm [13] for graph generation. During the forward diffusion process, GDSS injects Gaussian noise to both the node features and the adjacency matrix at each diffusion step. Then, a neural network is trained to learn the score function by minimizing the score-matching objective, which enables a reversion of graph data from noise via a reversed time diffusion process. However, such a directly borrowed diffusion model is incompatible with graph topology generation: unlike images which are feature-rich, the graph adjacency matrix is generally sparse and low-rank. Hence, injecting isotropic Gaussian noise into the sparsely connected parts of the adjacency matrix severely harms the graph data distribution and makes it hard to be recovered from Gaussian noise.

3 PRELIMINARIES

3.1 Notations

In this paper, we denote the probability space of interest as $(\Omega, \mathcal{F}, \mathbb{P})$ and $(\mathcal{F}_t)_{t \in \mathbb{R}}$ be a filtration, i.e. a sequence of increasing sub- σ -algebra of \mathcal{F} . Without specification, we denote $(\mathbf{B}_t)_{t \in \mathbb{R}}$ as the d -dimensional standard Brownian motion on the filtered probabilistic space $(\Omega, \mathcal{F}, \mathbb{P}, (\mathcal{F}_t)_{t \in \mathbb{R}})$. The distribution, support set and expectation of a random variable \mathbf{z} are defined as $\text{law}(\mathbf{z})$, $\text{supp}(\mathbf{z})$ and $\mathbb{E}[\mathbf{z}]$. $\mathbf{y}|\mathbf{z}$ denotes the distribution of \mathbf{y} conditioned on \mathbf{z} . $\text{Unif}(A)$

denotes the uniform distribution on a set A . $\|\cdot\|$ denotes the standard Euclid norm. $\|\cdot\|_\infty$ and $\|\cdot\|_{\text{lip}}$ denotes the supremum norm and Lipschitz norm of a function. $\lfloor \cdot \rfloor$ denotes the flooring function.

3.2 Score-based Generative Diffusion Model

A generative model refers to a mapping $g_\theta : \mathbb{R}^d \mapsto \mathbb{R}^d$, which maps a simple known priori π to a complicated data distribution \mathcal{D} . Once the model g_θ is sufficiently trained on N i.i.d samples from \mathcal{D} , denoted by \mathcal{S} , it enables us to generate plausible instances from \mathcal{D} directly, by sampling from $g_\theta(\varepsilon)$, $\varepsilon \sim \pi$. Unlike conventional generative models, e.g. VAE and GAN, which treat \mathcal{D} as a unilateral transformation of π , diffusion models consider the bilateral relation between \mathcal{D} and π from the perspective of SDE. Given an SDE travelling from \mathcal{D} to π , the corresponding reversed time SDE enables us to backtrack from noisy priori to the distribution of interest.

Lemma 3.1 (Forward Diffusion and Reversed Time SDE [17]).

The Forward Diffusion refers to the following SDE

$$\mathbf{z}_0 \sim \mathcal{D}, d\mathbf{z}_t = \mathbf{f}(\mathbf{z}_t, t)dt + \sigma_t d\mathbf{B}_t, t \in [0, 1], \quad (1)$$

where $\mathbf{f}(\cdot, t) : \mathbb{R}^d \mapsto \mathbb{R}^d$ is the drift function, $\sigma_t : [0, 1] \mapsto \mathbb{R}$ be a scalar diffusion function. Let $p_t(\cdot)$ be the probability density function of \mathbf{z}_t , then the Reversed Time SDE is given by

$$d\bar{\mathbf{z}}_t = (\mathbf{f}(\bar{\mathbf{z}}_t, t) - \sigma_t^2 \nabla \log p_t(\bar{\mathbf{z}}_t))d\bar{t} + \sigma_t d\bar{\mathbf{B}}_t, \quad (2)$$

$$\bar{\mathbf{z}}_1 \sim \mathbf{z}_1, t \in [0, 1],$$

where $d\bar{t} = -dt$ is the negative infinitesimal time step, $(\bar{\mathbf{B}}_t)_{t \in \mathbb{R}}$ is a reversed time Brownian motion w.r.t $(\Omega, \mathcal{F}, \mathbb{P}, (\bar{\mathcal{F}}_t)_{t \in \mathbb{R}})$, and $(\bar{\mathcal{F}}_t)_{t \in \mathbb{R}}$ is the corresponding decreasing filtration; $\nabla \log p_t(\cdot)$ is the score function.

During the forward diffusion process, with a carefully designed $\mathbf{f}(\cdot, t)$, the original data is perturbed by Gaussian noise with increasing magnitude, and it is assumed to be gradually corrupted to a truly noisy signal (priori), i.e. $\text{law}(\mathbf{z}_1) = \pi$. In order to draw new data from π via the reversed time SDE, one needs to learn the unknown score function $\nabla \log p_t(\cdot)$ with a neural network $s_\theta(\cdot) : \mathbb{R}^d \mapsto \mathbb{R}^d$, by minimizing the following explicit score matching error:

$$\mathcal{E}(\theta) \triangleq \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} \mathbb{E}_{\mathbf{z}_t | \mathbf{z}} \|s_\theta(\mathbf{z}_t) - \nabla \log p_t(\mathbf{z}_t)\|^2. \quad (3)$$

In practice, we employ a Gaussian priori π . For each sample $\mathbf{z}_j^i \in \mathcal{S}$, we first generate a sequence of corrupted data $\{\mathbf{z}_{t_j}^i\}_{j=1}^T$ by discretizing (1). To learn the score network $s_\theta(\cdot)$, one can minimize a more tractable denoising score matching objective $\hat{\mathcal{E}}(\theta)$

$$\hat{\mathcal{E}}(\theta) \triangleq \mathbb{E}_{\mathbf{z} \sim \text{Unif}(\mathcal{S})} \mathbb{E}_{\mathbf{z}_t | \mathbf{z}} \|s_\theta(\mathbf{z}_t) - \nabla \log p(\mathbf{z}_t | \mathbf{z})\|^2, \quad (4)$$

which has been proven to be equivalent to $\mathcal{E}(\theta)$ in [18]. Given a well trained score network $s_{\theta^*}(\cdot)$, one is able to generate plausible data from π via solving the learned reversed-time SDE

$$d\hat{\mathbf{z}}_t = (\mathbf{f}(\hat{\mathbf{z}}_t, t) - \sigma_t^2 s_{\theta^*}(\hat{\mathbf{z}}_t))d\bar{t} + \sigma_t d\bar{\mathbf{B}}_t, \quad (5)$$

$$\hat{\mathbf{z}}_1 \sim \pi, t \in [0, 1].$$

Ideally, the learned reversed time SDE should lead us towards \mathcal{D} , i.e. $\text{law}(\hat{\mathbf{z}}_0) = \mathcal{D}$.

4 METHODOLOGY

In this section, we establish the Graph Spectral Diffusion Model (GSDM) for fast and effective graph data generation. In Section 4.1, we briefly review the standard score-based graph diffusion model [10]. In Section 4.2, we formally introduce our GSDM algorithm and its α -quantile variants. In Section 4.3, we provide theoretical analyses to justify the efficacy of GSDM on graph data generation.

4.1 Standard Graph Diffusion Model

A graph with n nodes is defined as $\mathbf{G} \triangleq (\mathbf{X}, \mathbf{A}) \in \mathbb{R}^{n \times d} \times \mathbb{R}^{n \times n}$, where $\mathbf{X} \in \mathbb{R}^{n \times d}$ is the node feature matrix with dimension d and $\mathbf{A} \in \mathbb{R}^{n \times n}$ denotes the adjacency matrix. A graph generative model aims to learn the underlying data distribution, say $\mathbf{G} \sim \mathcal{G}$, which is a joint distribution of both \mathbf{X} and \mathbf{A} . Note that, if (\mathbf{X}, \mathbf{A}) is treated as a whole and omit the intrinsic graph structure, the aforementioned score-based generation framework can be parallelly extended to the graph generation setting, which yields the standard graph diffusion model, i.e. GDSS [10]. Roughly speaking, for each graph sample (\mathbf{X}, \mathbf{A}) , we first generate a sequence of perturbed graphs $\{(\mathbf{X}_{t_i}, \mathbf{A}_{t_i})\}_{i=1}^T$ via forward diffusion. Then, we train two score networks $s_\theta(\cdot)$ and $s_\phi(\cdot)$ to learn the score functions for both \mathbf{X}_t and \mathbf{A}_t , with which we can generate new data from π by running the reversed time SDE. Graph diffusion SDEs are defined as follows.

Definition 1 (Graph Diffusion SDEs with Disentangled Drift).
The Forward Graph Diffusion refers to the following SDE system

$$\begin{cases} d\mathbf{X}_t = \mathbf{f}^X(\mathbf{X}_t, t)dt + \sigma_{X,t}dB_t^X, \\ d\mathbf{A}_t = \mathbf{f}^A(\mathbf{A}_t, t)dt + \sigma_{A,t}dB_t^A, \end{cases} \quad (6)$$

$$(\mathbf{X}_0, \mathbf{A}_0) \sim \mathcal{G}, \quad t \in [0, 1],$$

where $\mathbf{f}^X(\cdot, t) : \mathbb{R}^{n \times d} \mapsto \mathbb{R}^{n \times d}$ and $\mathbf{f}^A(\cdot, t) : \mathbb{R}^{n \times n} \mapsto \mathbb{R}^{n \times n}$ is the drift functions for nodes feature and adjacency matrix; $\sigma_{X,t}, \sigma_{A,t}$ are the scalar diffusion terms (a.k.a noise schedule function); $(B_t^X)_{t \in \mathbb{R}}$ and $(B_t^A)_{t \in \mathbb{R}}$ are standard Brownian motions on $\mathbb{R}^{n \times d}$ and $\mathbb{R}^{n \times n}$, respectively.

To alleviate the computational burden of calculating the drift w.r.t the high dimensional \mathbf{G} , the drift term of forward graph diffusion is disentangled into $\mathbf{f}^X(\cdot, t)$ and $\mathbf{f}^A(\cdot, t)$. Again, Lemma 3.1 guarantees the existence of the reversed time SDEs for graph diffusion.

Corollary 1 (Reversed Time SDEs for Graph Diffusion). The reversed time SDE system of (6) is given by

$$\begin{cases} d\bar{\mathbf{X}}_t = (\mathbf{f}^X(\bar{\mathbf{X}}_t, t) - \sigma_{X,t}^2 \nabla_{\mathbf{X}} \log p_t(\bar{\mathbf{X}}_t, \bar{\mathbf{A}}_t)) d\bar{t} + \sigma_{X,t} dB_t^X, \\ d\bar{\mathbf{A}}_t = (\mathbf{f}^A(\bar{\mathbf{A}}_t, t) - \sigma_{A,t}^2 \nabla_{\mathbf{A}} \log p_t(\bar{\mathbf{X}}_t, \bar{\mathbf{A}}_t)) d\bar{t} + \sigma_{A,t} dB_t^A, \end{cases} \quad (7)$$

$$(\bar{\mathbf{X}}_0, \bar{\mathbf{A}}_1) \sim \pi, \quad t \in [0, 1],$$

where $d\bar{t} = -dt$ is the negative infinitesimal time step; $(\bar{B}_t^X)_{t \in \mathbb{R}}$ and $(\bar{B}_t^A)_{t \in \mathbb{R}}$ are the reversed time standard Brownian motions induced by (6).

The disentanglement of the drift functions implies the conditional independence $\mathbf{X}_t \perp \mathbf{A}_t | \mathbf{G}_0$, with which we can decompose $p_t(\bar{\mathbf{X}}_t, \bar{\mathbf{A}}_t)$ into $p_{t|0}(\mathbf{X}_t | \mathbf{X}_0) \cdot p_{t|0}(\mathbf{A}_t | \mathbf{A}_0)$, where $p_{t|0}(\cdot)$ denotes the density function of $\mathbf{G}_t | \mathbf{G}_0$. As

proposed in [10], such conditional independence reduces the denoising score matching objective to a simpler form

$$\hat{\mathcal{E}}(\theta) \triangleq \mathbb{E}_{\mathbf{G} \sim \text{Unif}(\mathcal{S})} \mathbb{E}_{\mathbf{G}_t | \mathbf{G}} \|s_\theta(\mathbf{G}_t) - \nabla \log p_{t|0}(\mathbf{X}_t | \mathbf{X}_0)\|^2, \quad (8)$$

$$\hat{\mathcal{E}}(\phi) \triangleq \mathbb{E}_{\mathbf{G} \sim \text{Unif}(\mathcal{S})} \mathbb{E}_{\mathbf{G}_t | \mathbf{G}} \|s_\phi(\mathbf{G}_t) - \nabla \log p_{t|0}(\mathbf{A}_t | \mathbf{A}_0)\|^2. \quad (9)$$

Hence, the training and sampling procedures can be directly borrowed from standard score-based models.

While GDSS is the first attempt at leveraging diffusion models on graph generation, its performance is hindered by the brute force application of diffusion. For sparsely connected graphs, while the distribution of node features varies across datasets, the distribution of graph topology, i.e. adjacency matrix, resides in a low dimensional manifold. As mentioned in Section 1, an evident pattern of the adjacency matrices also implies that the true distribution of \mathbf{A} is of low rank. In this case, the score matching objective fails to provide consistency estimators. Although running a full rank diffusion on $\mathbf{A} \in \mathbb{R}^{n \times n}$ alleviates this issue by extending the support of corrupted data from the manifold to the full space, it inevitably introduces lethal noise to regions of zero probability density. As a consequence, the signal-to-noise ratio of regions out of $\text{supp}(\mathbf{A})$ is essentially zero, which is a catastrophe for training the denoising score network. Note that such full-rank diffusion is also inappropriate for the densely connected graph generation. This is because an isotropically corrupted adjacency matrix encourages delusive message passing on sparsely connected parts of the graph, which is destructive to the graph message passing pattern. Thus, standard diffusion can severely impair representation learning for sparse graph regions.

4.2 Graph Spectral Diffusion Model

To address these notorious yet ubiquitous issues, we novelly propose the Graph Spectral Diffusion Model. For graph topology generation, in contrast to GDSS which is driven by a full-rank diffusion on the whole space $\mathbb{R}^{n \times n}$, our GSDM leverages low-rank diffusion SDEs on the n -dimensional spectrum manifold, e.g. the span of n -eigenvalues of \mathbf{A} . As we shall see later, GSDM achieves both robustness and computational efficiency, by exploiting the graph spectrum structure and running diffusion on an information-concentrated manifold.

Definition 2 (Graph Spectral Diffusion SDEs).

Let the spectral decomposition of \mathbf{A} be $\mathbf{U}\Lambda\mathbf{U}^\top$, where columns of \mathbf{U} are the orthonormal eigenvectors and Λ be the diagonal eigenvalue matrix, i.e. spectrum. The Forward Spectral Diffusion refers to the following SDE system

$$\begin{cases} d\mathbf{X}_t = \mathbf{f}^X(\mathbf{X}_t, t)dt + \sigma_{X,t}dB_t^X, \\ d\mathbf{A}_t = \mathbf{f}^A(\mathbf{A}_t, t)dt + \sigma_{A,t}dB_t^A, \end{cases} \quad (10)$$

$$(\mathbf{X}_0, \mathbf{A}_0) \sim \mathcal{G}, \quad \mathbf{A}_0 = \mathbf{U}_0 \Lambda_0 \mathbf{U}_0^\top, \quad t \in [0, 1],$$

where $\mathbf{f}^X(\cdot, t) : \mathbb{R}^{n \times d} \mapsto \mathbb{R}^{n \times d}$ is the drift for nodes feature; $\mathbf{f}^A(\cdot, t) : \mathbb{R}^n \mapsto \mathbb{R}^n$ is the drift for spectrum, which only acts on the diagonal entries; $\sigma_{X,t}, \sigma_{A,t}$ are the scalar diffusion terms (a.k.a noise schedule functions); $(B_t^X)_{t \in \mathbb{R}}$ and $(B_t^A)_{t \in \mathbb{R}}$ are standard Brownian motions on $\mathbb{R}^{n \times d}$ and \mathbb{R}^n , respectively; and $\mathbf{W}_t^\Lambda \triangleq \text{diag}(B_t^A)$ is a diagonal Brownian motion.

By comparing (10) with (6), one should notice that the initial eigenvector matrix \mathbf{U}_0 is fixed along the spectral diffusion. As will be illustrated later, the evolution of $\mathbf{A}_t \triangleq \mathbf{U}_0 \mathbf{A}_t \mathbf{U}_0^\top$ is driven by a n -dimensional Gaussian process. Hence, we prevent the corrupted adjacency matrix from rampaging around the full space. We are now ready to establish the reversed time spectral diffusion SDEs.

Corollary 2 (Reversed Time Spectral Diffusion SDEs). *The reversed time Spectral Diffusion SDE system of (10) is given by*

$$\begin{cases} d\bar{\mathbf{X}}_t = (\mathbf{f}^X(\bar{\mathbf{X}}_t, t) - \sigma_{X,t}^2 \nabla_{\mathbf{X}} \log p_t(\bar{\mathbf{X}}_t, \bar{\Lambda}_t)) dt + \sigma_{X,t} d\bar{\mathbf{B}}_t^X, \\ d\bar{\Lambda}_t = (\mathbf{f}^\Lambda(\bar{\Lambda}_t, t) - \sigma_{\Lambda,t}^2 \nabla_{\Lambda} \log p_t(\bar{\mathbf{X}}_t, \bar{\Lambda}_t)) dt + \sigma_{\Lambda,t} d\bar{\mathbf{W}}_t^\Lambda, \end{cases} \quad (11)$$

$(\bar{\mathbf{X}}_1, \bar{\Lambda}_1) \sim \pi, t \in [0, 1],$

where $d\bar{t} = -dt$ is the negative infinitesimal time step; $(\bar{\mathbf{B}}_t^X)_{t \in \mathbb{R}}$ and $(\bar{\mathbf{B}}_t^\Lambda)_{t \in \mathbb{R}}$ are reversed time standard Brownian motions induced by (6); and $\bar{\mathbf{W}}_t^\Lambda \triangleq \text{diag}(\bar{\mathbf{B}}_t^\Lambda)$.

Since \mathbf{U} is no longer involved in (11), the boundary condition is only imposed on the joint distribution of $(\mathbf{X}_1, \Lambda_1)$ such that $\text{law}(\mathbf{X}_1, \Lambda_1) = \pi$. This assumption implies that, the authentic distribution $(\mathbf{X}_0, \Lambda_0)$ can be recovered from a priori π by the reversed time spectral diffusion SDE. According to score matching techniques [18], we can train two score networks $s_\theta(\cdot, \cdot), s_\phi(\cdot, \cdot)$ to learn the score functions $\nabla_{\mathbf{X}} \log p_t(\cdot, \cdot), \nabla_{\Lambda} \log p_t(\cdot, \cdot)$ via minimizing

$$\begin{aligned} \hat{\mathcal{E}}(\theta) &\triangleq \mathbb{E}_{\mathbf{G} \sim \text{Unif}(\mathcal{S})} \mathbb{E}_{\mathbf{X}_t | \mathbf{G}} \|s_\theta(\mathbf{X}_t, \Lambda_t) - \nabla \log p_{t|0}(\mathbf{X}_t | \mathbf{X}_0)\|^2, \\ \hat{\mathcal{E}}(\phi) &\triangleq \mathbb{E}_{\mathbf{G} \sim \text{Unif}(\mathcal{S})} \mathbb{E}_{\Lambda_t | \mathbf{G}} \|s_\phi(\mathbf{X}_t, \Lambda_t) - \nabla \log p_{t|0}(\Lambda_t | \Lambda_0)\|^2. \end{aligned}$$

In nutshell, the proposed GSDM is summarized as three main steps. Details of GSDM can be found Algorithm 1 and Algorithm 2 in Appendix B.1.

1. Run forward spectral diffusion model on (\mathbf{X}, Λ) by (10). Train two score networks $s_\theta(\cdot), s_\phi(\cdot)$ to learn score functions shown up in (11).
2. Generate plausible $(\hat{\mathbf{X}}_0, \hat{\Lambda}_0)$ from π , via reversing the spectral diffusion SDEs from $t = 1$ to $t = 0$, with estimated score functions $s_\theta(\hat{\mathbf{X}}_t, \hat{\Lambda}_t), s_\phi(\hat{\mathbf{X}}_t, \hat{\Lambda}_t)$.
3. Generate plausible adjacency matrix via $\hat{\mathbf{A}} = \hat{\mathbf{U}} \hat{\Lambda}_0 \hat{\mathbf{U}}^\top$, where $\hat{\mathbf{U}}$ is uniformly sampled from the observed eigenvector matrices.

Since the full adjacency matrices are not involved in the computation of diffusion SDEs, GSDM achieves significant acceleration in both training and sampling. Moreover, one can further enhance the computational efficiency, by confining the spectral diffusion to the top- k largest eigenvalues of \mathbf{A} , where $k \triangleq [\alpha n]$. Suppose Λ^k is the truncated eigenvalue matrix, where only the top- k diagonal entries of Λ are preserved, we can define the α -quantile GSDM by substituting the occurrence of Λ in GSDM with $\Lambda^{(k)}$. As will be seen in the following section, α -quantile GSDM exhibits evidently faster processing speed and comparable performance to GSDM.

4.3 Theoretical Analysis

Here, we provide supportive theoretical evidence for the efficacy of GSDM. In Proposition 1, we first study the low-rank structure of the spectral diffusion SDE of adjacency matrices. In Proposition 2, we further prove that our proposed spectral diffusion enjoys a sharper reconstruction error bound than the standard graph diffusion model.

Proposition 1 (Spectral Diffusion SDEs on Adjacency Matrix). *Suppose $\mathbf{f}_\Lambda(\Lambda, t) \triangleq -\sigma_{\Lambda,t}^2/2\Lambda$. The spectral diffusion SDE system (10) induces an n -dimensional SDE system of the adjacency matrix \mathbf{A} on the full space $\mathbb{R}^{n \times n}$. Following previous notations, the forward diffusion SDE is given by*

$$\begin{cases} d\mathbf{X}_t = \mathbf{f}^X(\mathbf{X}_t, t) dt + \sigma_{X,t} d\mathbf{B}_t^X, \\ d\Lambda_t = -\frac{1}{2}\sigma_{\Lambda,t}^2 \Lambda_t dt + \sigma_{\Lambda,t} d\mathbf{M}_t, \end{cases} \quad (12)$$

where $(\mathbf{M}_t)_{t \in [0,1]}$ is a n -dimensional centered Gaussian process on $\mathbb{R}^{n \times n}$, with zero mean and covariance kernel $\mathcal{K}(s, t) : [0, 1] \times [0, 1] \mapsto \mathbb{R}^{n \times n \times n \times n}$ as

$$\mathcal{K}(s, t)_{i,j,k,l} = \min(s, t) \cdot \sum_{h=1}^n \mathbf{U}_0[i, h] \mathbf{U}_0[j, h] \mathbf{U}_0[k, h] \mathbf{U}_0[l, h].$$

Hence, the conditional distribution of \mathbf{A}_t on \mathbf{A}_0 is Gaussian

$$\mathbf{A}_t | \mathbf{A}_0 \sim N(\mathbf{A}_t; \mathbf{A}_0 e^{-\frac{1}{2} \int_0^t \sigma_\tau^2 d\tau}, (1 - e^{-\int_0^t \sigma_\tau^2 d\tau}) \mathcal{K}(1, 1)), \quad (13)$$

which admits a closed-form probability density function.

Remark 1. *The proof is postponed to Appendix A.1. Proposition 1 shows that our spectral diffusion framework substantially recasts the evolution of adjacency matrix, by driving the n^2 -dimensional SDE with n -dimensional noise $(\mathbf{M}_t)_{t \in \mathbb{R}}$. Guided by the graph spectrum, the diffusion is concentrated to the salient parts of $\text{supp}(\mathbf{A})$ to prevent introducing irreducible noise to the out-of-support regions.*

The central question for graph generation is how to measure the quality of the synthesis data that is recovered from noise with the learned score function. The key to this question is to establish the reconstruction error bound, i.e. the expected error between the data reconstructed with ground truth score $\nabla \log p_t(\cdot)$ and the learned scores $s_\phi(\cdot)$. For graph topology generation, our GSDM is proven to enjoy a sharper reconstruction bound than standard graph diffusion. Detailed proof can be found in Appendix A.2.

Proposition 2 (Reconstruction Bounds for Adjacency Matrix Generation). *Following the previous notations, we define the estimated reversed time SDEs for standard and spectral graph diffusion models as*

$$\begin{aligned} d\hat{\mathbf{A}}_t^{\text{full}} &= \left(-\frac{1}{2}\sigma_t^2 \hat{\mathbf{A}}_t^{\text{full}} - \sigma_t^2 s_\phi(\hat{\mathbf{A}}_t^{\text{full}}, t) \right) dt + \sigma_t d\bar{\mathbf{B}}_t^A, \\ d\hat{\mathbf{A}}_t^{\text{spec}} &= \left(-\frac{1}{2}\sigma_t^2 \hat{\mathbf{A}}_t^{\text{spec}} - \sigma_t^2 s_\varphi(\hat{\mathbf{A}}_t^{\text{spec}}, t) \right) dt + \sigma_t d\bar{\mathbf{M}}_t, \end{aligned}$$

where both generation methods share the same drift term and noise schedule $(\sigma_t)_{t \in [0,1]}$. We further assume that $\|s_\phi(\cdot, t)\|_{\text{lip}} = \mathcal{O}(\mathbb{E}_{|\mathbf{A}_0|} \|\nabla_{\Lambda} \log p_{t|0}(\hat{\mathbf{A}}_t^{\text{full}})\|_{\text{lip}})$ and $\|s_\varphi(\cdot, t)\|_{\text{lip}} = \mathcal{O}(\mathbb{E}_{|\mathbf{A}_0|} \|\nabla_{\Lambda} \log p_{t|0}(\hat{\mathbf{A}}_t^{\text{spec}})\|_{\text{lip}})$ holds almost surely. By reversing both SDEs from $t = 1$ to $t = 0$, we obtain $\hat{\mathbf{A}}_t^{\text{full}}$ and $\hat{\mathbf{A}}_t^{\text{spec}}$, which are two reconstructions of the authentic \mathbf{A}_0 , with expected reconstruction errors bounded by

$$\begin{aligned} &\mathbb{E} \|\mathbf{A}_0 - \hat{\mathbf{A}}_0^{\text{full}}\|^2 \\ &\leq M \mathcal{E}(\phi) \cdot \left(1 + n^2 K \int_0^1 \Sigma_t^{-2} \exp \left(n^2 K \int_t^1 \Sigma_s^{-2} ds \right) dt \right), \\ &\mathbb{E} \|\mathbf{A}_0 - \hat{\mathbf{A}}_0^{\text{spec}}\|^2 \\ &\leq M \mathcal{E}(\varphi) \cdot \left(1 + n K \int_0^1 \Sigma_t^{-2} \exp \left(n K \int_t^1 \Sigma_s^{-2} ds \right) dt \right), \end{aligned}$$

where $M \triangleq C^2 \|\sigma\|_\infty^4$; $K \triangleq 2ML/\mathbb{E}\|\mathbf{A}_0\|_{2,2}$; C, L are absolute constants; $\Sigma_t^2 \triangleq 1 - e^{-\int_0^t \sigma_s^2 ds}$; $\mathcal{E}(\cdot)$ is the expected score matching objective defined in (3).

Remark 2. While the error bound of $\hat{\mathbf{A}}_0^{\text{full}}$ is at the order of $\mathcal{O}(n^2 \exp(n^2))$, $\hat{\mathbf{A}}_0^{\text{spec}}$ exhibits a much sharper bound of order $\mathcal{O}(n \exp(n))$. For large-scale graphs with a large number of nodes n , our proposed GSDM enjoys a substantially better performance guarantee, which coincides with our numerical results. Moreover, under additional conditions, the error bounds can be significantly improved for well-trained score networks ϕ^*, φ^* , by showing that $\mathcal{E}(\phi^*) \ll \mathcal{E}(\phi)$ as in Proposition 3. The detailed proof can be found in Appendix A.3.

Definition 3 (β -smooth). $f : \mathbb{R}^m \mapsto \mathbb{R}^n$ is called β -smooth if and only if

$$\|f(x_1) - f(x_2) - \nabla f(x_2)^\top (x_1 - x_2)\| \leq \frac{\beta}{2} \|x_1 - x_2\|^2$$

holds for $\forall x_1, x_2 \in \mathbb{R}^m$.

Proposition 3 (Convergence of Score Matching Objective Minimization). Recall that the score matching objective of a generic sample $\mathbf{Z} \in \mathbb{R}^d$ ($\hat{\mathbf{A}}^{\text{full}}$ or $\hat{\mathbf{A}}^{\text{spec}}$) is defined as

$$\mathcal{E}(\theta; \mathbf{Z}) \triangleq \mathbb{E}_{\mathbf{Z}_t | \mathbf{Z}} \|s_\theta(\mathbf{Z}_t, t) - \nabla_{\mathbf{Z}} \log p_t(\mathbf{Z}_t)\|^2, \quad (14)$$

and the expected and empirical score matching error are defined as

$$\mathcal{E}(\theta) \triangleq \mathbb{E}_{\mathbf{Z} \sim \mathcal{D}} \mathcal{E}(\theta; \mathbf{Z}) \quad (15)$$

$$\widehat{\mathcal{E}}(\theta; \mathcal{S}_N) \triangleq \mathbb{E}_{\mathbf{Z} \sim \mathcal{S}_N} \mathcal{E}(\theta; \mathbf{Z}) \quad (16)$$

where \mathcal{D} is the population distribution of \mathbf{Z} , and $\mathcal{S} \triangleq \{\mathbf{Z}^i\}_{i=1}^N$ is the uniform distribution over an i.i.d sampled training dataset of size N . Suppose $\mathcal{E}(\theta)$ is minimized via running standard Stochastic Gradient Descent (SGD) on training data, i.e. at the k -th iteration, θ_k is updated on a mini-batch of size b

$$\theta_{k+1} \triangleq \theta_k - \eta \nabla_\theta \widehat{\mathcal{E}}(\theta; \mathcal{S}_b). \quad (17)$$

Assume that almost surely (w.r.t \mathbf{Z}), $\mathcal{E}(\cdot; \mathbf{Z})$ is β -smooth, and the tangent kernel $\mathbf{K}_\theta(\mathcal{S}) \in \mathbb{R}^{Nd \times Nd}$ of $s_\theta(\cdot)$ satisfies

$$\lambda_{\min}(\mathbf{K}_\theta(\mathcal{S})) \geq \lambda > 0, \quad \theta \in B(\theta_0, R),$$

$$\mathbf{K}_\theta(\mathcal{S})[Ni + 1 : N(i + 1), Nj + 1 : N(j + 1)]$$

$$\triangleq \nabla_\theta s_\theta(\mathbf{Z}^i)^\top \nabla_\theta s_\theta(\mathbf{Z}^j),$$

with $R = 2N\sqrt{2\beta\mathcal{E}(\theta_0)}/(\mu\delta)$, $\delta > 0$. Then, with probability $1 - \delta$ over the choice of mini-batch \mathcal{S}_b , SGD with a learning rate $\eta \leq \frac{\lambda/N}{N\beta(N^2\beta + \lambda(b-1)/N)}$ converges to a global solution in the ball $B(\theta_0, R)$ with exponential convergence rate

$$\mathcal{E}(\theta_k) \leq \left(1 - \frac{\lambda b \eta}{N}\right)^k \mathcal{E}(\theta_0). \quad (18)$$

5 EXPERIMENTS

In this section, we evaluate our proposed GDSM with state-of-the-art graph generative baselines on two types of graph generation tasks: generic graph generation and molecule generation, over several benchmark datasets. We also conduct extensive ablation studies and visualizations to further illustrate the effectiveness and efficiency of GDSM.

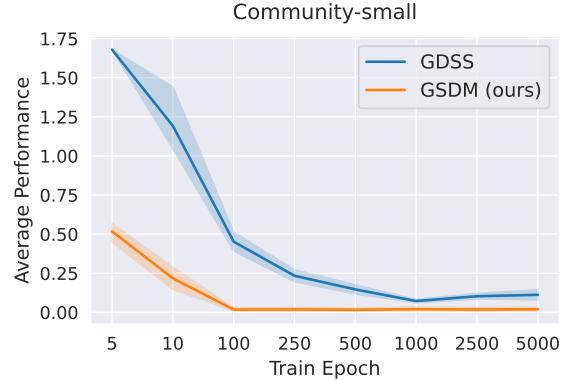


Fig. 3: GSDM enjoys a significantly faster convergence rate than GDSS. On the community-small dataset, our GSDM reaches the SOTA performance within 100 training epochs.

5.1 Generic Graph Generation 数据集

Datasets We test our model on three generic datasets with different scales, and we use N to denote the number of nodes: (1) **Community-small** ($12 \leq N \leq 20$): contains 100 small community graphs. (2) **Enzymes** ($10 \leq N \leq 125$): contains 578 protein graphs which represent the protein tertiary structures of the enzymes from the BRENDA database. (3) **Grid** ($100 \leq N \leq 400$): contains 100 standard 2D grid graphs. To fairly compare our model with baselines, we adopt the experimental and evaluation setting from [9] with the same train/test split.

Baselines We compare our model with well-known or state-of-the-art graph generation methods, which can be categorized into auto-regressive models and one-shot models. The **auto-regressive model** refers to sequential generation, which constructs a graph via a set of consecutive steps, usually done nodes by nodes and edges by edges [14]. Under this category, we include DeepGMG [19], GraphRNN [9], GraphAF [4], and GraphDF [20]. The **one-shot model** refers to building a probabilistic graph model that can generate all nodes and edges in one shot [14]. Under this category, we include GraphVAE [8], Graph Normalizing Flow(GNF) [21], EDP-GNN [12], and GDSS [10].

Metrics We adopt the **maximum mean discrepancy (MMD)** to compare the distributions of graph statistics, such as the degree, the clustering coefficient, and the number of occurrences of orbits with 4 nodes [9], [10] between the same number of generated and test graphs.

Results We show the results in Table 1. The results show that our proposed GSDM significantly outperforms both the auto-regressive baselines and one-shot baseline methods. Specifically, GDSS is the SOTA graph diffusion model which performs the diffusion in the whole space of the graph data. Our method substantially outperforms GDSS in terms of both average performance and convergence rate (Figure 3), which demonstrates the advantages of performing SDEs in the spectrum of the graph compared to the whole space.

5.2 Molecules Generation

Besides generic graph generation, our model can also generate organic molecules through our proposed reverse diffusion process. We test our model with two well-known

TABLE 1: **Generation results on the generic graph datasets.** We report the MMD distances between the test datasets and generated graphs. The best results are highlighted in bold (the smaller the better). Hyphen (-) denotes out-of-resources that take more than 10 days or are not applicable due to memory issues.

		Community-small				Enzymes				Grid			
		Synthetic, $12 \leq V \leq 20$				Real, $10 \leq V \leq 125$				Synthetic, $100 \leq V \leq 400$			
		Deg. \downarrow	Clus. \downarrow	Orbit \downarrow	Avg. \downarrow	Deg. \downarrow	Clus. \downarrow	Orbit \downarrow	Avg. \downarrow	Deg. \downarrow	Clus. \downarrow	Orbit \downarrow	Avg. \downarrow
Autoreg.	DeepGMG [19]	0.220	0.950	0.400	0.523	-	-	-	-	-	-	-	-
	GraphRNN [9]	0.080	0.120	0.040	0.080	0.017	0.043	0.021	0.043	-	-	-	-
	GraphAF [4]	0.18	0.20	0.02	0.133	1.669	1.283	0.266	1.073	-	-	-	-
One-shot	GraphDF [20]	0.06	0.12	0.03	0.070	1.503	1.061	0.202	0.922	-	-	-	-
	GraphVAE [8]	0.350	0.980	0.540	0.623	1.369	0.629	0.191	0.730	1.619	0.0	0.919	0.846
	GNF [21]	0.200	0.200	0.110	0.170	-	-	-	-	-	-	-	-
	EDP-GNN [12]	0.053	0.144	0.026	0.074	0.023	0.268	0.082	0.124	0.455	0.238	0.328	0.340
	GDSS ¹ [10]	0.045	0.086	0.007	0.046	0.026	0.102	0.009	0.046	0.111	0.005	0.070	0.062
Ours		0.011	0.015	0.001	0.009	0.013	0.088	0.01	0.037	0.002	0.0	0.0	0.0007

¹ The average results of the Enzymes dataset reported in the GDSS original paper is 0.032. However, the best result we can obtain using the author’s released code and checkpoint with careful fine-tuning is 0.046.

TABLE 2: **Generation results on the QM9 and ZINC250k datasets.** Results are the means of three different runs, and the best results are highlighted in bold. Values denoted by * are taken from the respective original papers. Other results are obtained by running open-source codes. Val. w/o corr. denotes the Validity w/o correction metric, and values that do not exceed 50% are underlined.

		QM9					ZINC250k				
Method		Validity (%) \uparrow	Val. w/o corr. (%) \uparrow	NSPDK \downarrow	FCD \downarrow	Time (s) \downarrow	Validity (%) \uparrow	Val. w/o corr. (%) \uparrow	NSPDK \downarrow	FCD \downarrow	Time (s) \downarrow
Autoreg.	GraphAF [4]	100	<u>67*</u>	0.020	5.268	<u>2.28e³</u>	100	<u>68*</u>	0.044	16.289	<u>5.72e³</u>
	GraphAF+FC	100	<u>74.43</u>	0.021	5.625	<u>2.32e³</u>	100	<u>68.47</u>	0.044	16.023	<u>5.91e³</u>
	GraphDF [20]	100	<u>82.67*</u>	0.063	10.816	<u>5.08e⁴</u>	100	<u>89.03*</u>	0.176	34.202	<u>5.87e⁴</u>
One-shot	GraphDF+FC	100	93.88	0.064	10.928	<u>4.72e⁴</u>	100	90.61	0.177	33.546	<u>5.79e⁴</u>
	MoFlow [3]	100	91.36	0.017	4.467	4.58	100	63.11	0.046	20.931	25.9
	EDP-GNN [12]	100	<u>47.52</u>	0.005	2.680	<u>4.13e³</u>	100	82.97	0.049	16.737	<u>8.41e³</u>
	GraphEBM [22]	100	<u>8.22</u>	0.030	6.143	35.33	100	<u>5.29</u>	0.212	35.471	<u>53.72</u>
	GDSS [10]	100	95.72*	<u>0.003*</u>	2.900*	<u>1.06e²</u>	100	<u>97.01*</u>	0.019*	14.656*	2.11e ³
Ours		100	99.9	0.003	2.650	18.02	100	92.70	0.017	12.956	45.91

数据集

molecule datasets: QM9 [23] and ZINC250K [24]. Following previous works [10], [20], the molecules are kekulized by the RDKit library [25] with hydrogen atoms removed. We evaluate the quality of 10,000 generated molecules with Frechet ChemNet Distance (FCD) [26], Neighborhood subgraph pairwise distance kernel (NSPDK) MMD [27], validity w/o correction, and the generation time. FCD computes the distance between the testing and the generated molecules using the activations of the penultimate layer of the ChemNet. (NSPDK) MMD computes the MMD between the generated and the testing set which takes into account both the node and edge features for evaluation. Generally speaking, FCD measures the generation quality in the view of molecules in the chemical space, while NSPDK MMD evaluates the generation quality from the graph structure perspective. Besides, following [10], we also include the validity w/o correction as another metric to explicitly evaluate the quality of molecule generation prior to the correction procedure. It computes the fraction of the number of valid molecules without valency correction or edge resampling over the total number of generated molecules. In contrast, validity measures the fraction of the valid molecules after the correction phase. Generation time measures the time for generating 10,000 molecules in the form of RDKit. It is a salient measure of the practicability of molecule generation, especially for macromolecule generation.

Baselines We compare our model with the state-of-the-art molecule generation models. The baselines include SOTA

auto-regressive models: GraphAF [4] is a flow-based model, and GraphDF [20] is a flow-based model using discrete latent variables. Following GDSS [10], we modify the architecture of GraphAF and GraphDF to consider formal charges in the molecule generation, denoted as GraphAF+FC and GraphDF+FC, for fair comparisons. For the one-shot model, we include MoFlow [3], which is a flow-based model; EDP-GNN [12] and GDSS [12] which are both diffusion models.

Results We show the results in Table 2. Evidently, GSDM achieves the highest performance under most of the metrics. The highest scores in NSPDK and FCD show that GSDM is able to generate molecules that have close data distributions to the real molecules in both the chemical space and graph space. Especially, our model outperforms GDSS, in most of the metrics, verifying that our proposed spectral diffusion is not only suitable for generic graph generation but also advisable for molecule designs. Apart from the visualization of generated graphs in Figure 2, visualizations of generated molecules are included in Appendix B.3, showing that the proposed GSDM is capable of generating generic graphs as well as molecules of high quality.

Time Complexity One of the main advantages of our GSDM compared to other diffusion models (e.g. EDP-GNN and GDSS) is the efficiency in generating molecules. We show the time spent (in seconds) for generating 10,000 molecules in Table 2, demonstrating that our GSDM takes a significantly lower time in the inference process compared to EDP-GNN and GDSS. Especially, compared to

TABLE 3: **Ablation study on the α -quantile eigenvalues.** The metrics used are the same as in Table 1. The results that **surpass** the baseline methods in Table 1 are highlighted in bold (the smaller the better). Avg.% denotes the percentage of average scores achieved compared to using the whole eigenvalues and eigenvectors set.

α	Community-small					Enzymes					Grid				
	Synthetic, $12 \leq V \leq 20$					Real, $10 \leq V \leq 125$					Synthetic, $100 \leq V \leq 400$				
	Deg. \downarrow	Clus. \downarrow	Orbit \downarrow	Avg. \downarrow	Avg.% \uparrow	Deg. \downarrow	Clus. \downarrow	Orbit \downarrow	Avg. \downarrow	Avg.% \uparrow	Deg. \downarrow	Clus. \downarrow	Orbit \downarrow	Avg. \downarrow	Avg.% \uparrow
10%	1.010	1.105	0.227	0.781	0.012	0.446	0.608	0.059	0.371	0.029	1.996	0	1.013	1.003	0.001
20%	0.450	1.102	0.102	0.551	0.016	0.056	0.220	0.012	0.096	0.395	1.996	0	1.013	1.003	0.001
30%	0.085	0.421	0.030	0.179	0.050	0.011	0.093	0.012	0.039	0.982	0.107	0	0.044	0.050	0.012
40%	0.039	0.056	0.006	0.034	0.268	0.010	0.093	0.011	0.038	1.000	0.018	0	0.001	0.006	0.095
50%	0.022	0.024	0.002	0.016	0.563	0.010	0.093	0.011	0.038	1.000	0.007	0	0.001	0.002	0.225
60%	0.014	0.019	0.001	0.011	0.794	0.010	0.093	0.011	0.038	1.000	0.002	0	0	0.001	1.000
70%	0.012	0.019	0.001	0.011	0.844	0.010	0.093	0.011	0.038	1.000	0.002	0	0	0.001	1.000
80%	0.011	0.016	0.001	0.009	0.964	0.010	0.093	0.011	0.038	1.000	0.002	0	0	0.001	1.000
90%	0.011	0.015	0.001	0.009	1.000	0.010	0.093	0.011	0.038	1.000	0.002	0	0	0.001	1.000
100%	0.011	0.015	0.001	0.009	1.000	0.010	0.093	0.011	0.038	1.000	0.002	0	0	0.001	1.000

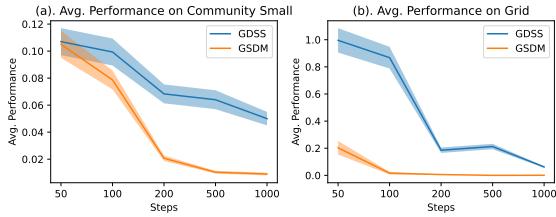


Fig. 4: Ablation studies on different diffusion step numbers.

GDSS which requires $1.06e^3$ and $2.11e^3$ seconds to generate 10,000 molecule graphs according to QM9 and Zinc250k, our model only takes 18.02 and 45.91 seconds, which are **58 \times** and **46 \times** faster respectively. This phenomenon verifies our methodology in designing the spectral diffusion that not only the spectral diffusion is more effective than the whole space diffusion, but also more efficient. Such an improvement is crucial for numerous applications such as drug design and material analysis.

5.3 Ablation Studies

To evaluate the effectiveness of GSDM, we conduct extensive ablation experiments and display elaborate results and discussions in this section.

Ablation studies on the number of diffusion steps. To analyse the robustness of our proposed GSDM during the reconstruction procedure (i.e. reverse diffusion process), we compare the performance of our model with GDSS on varying numbers of diffusion steps using the Community small and Grid dataset, which represents small and large generic graphs respectively. The default steps for our model and GDSS during training and evaluation are 1,000 steps. Ideally, for diffusion models, with the reduction in step numbers, the performance of the models decreases, while the model of higher robustness still exhibits a higher performance under different numbers of diffusion steps. The results are shown in Figure 4. We can observe that GSDM outperforms GDSS in both datasets under different diffusion steps. Specifically, in the Grid dataset, GSDM consistently exhibits significantly higher performance than GDSS. In the community small dataset, although GSDM and GDSS achieve similar performance at 50 steps, GSDM becomes much more accurate when the number of steps is higher than 100.

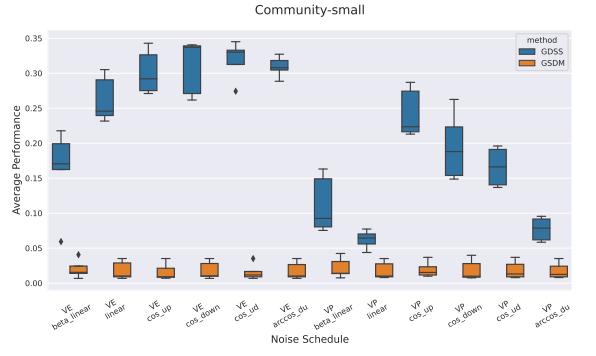


Fig. 5: Ablation studies on diffusion schedules. The y-axis denotes the average performance score (\downarrow) and each column denotes a diffusion schedule configuration. Each box plot summarizes the results of 5 random trials.

Ablation studies on the type of diffusion schedules. As mentioned in [13], [28] and [29], the performance of diffusion models highly relies on the diffusion schedule, i.e. the scheduled noise insertion process or discretization of certain types of SDEs, which lies at the heart of both training and inference phases. According to the taxonomy proposed in [13] and [11], mainstream artificial diffusion schedules are categorized to be either Variance Exploding (VE) or Variance Preserving (VP). To improve model robustness against the choice of diffusion schedules, [29] and [28] designed diffusion models with optimizable diffusion schedules. In this subsection, we conduct systematic experiments on several datasets to test the robustness of our proposed GSDM against different diffusion schedules. As illustrated in Figure 5, the state-of-the-art performance of our GSDM is immune to the choice of diffusion schedule and it frees us from tedious hyperparameter searching. When comparing with GDSS, our GSDM exhibits evidently better average performance with much smaller variance across various configurations. More results and experiment details are included in Appendix B.2.

Ablation studies on low-rank approximations of the score-function. As mentioned in Section 4.2, one can further accelerate the training and sampling phases, by confining diffusion to the partial spectrum of the graph adjacency matrix. Such variants are coined as α -quantile GSDM in Section 4.2.

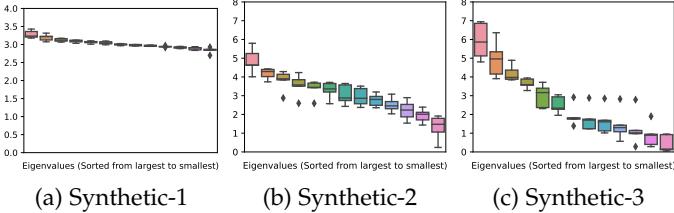


Fig. 6: The eigenvalue distribution of 3 synthetic datasets.

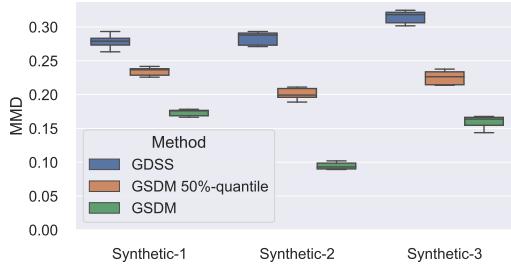


Fig. 7: Experiments on synthetic datasets of which the eigenvalues' distributions are shown in Figure 6. The y-axis denotes the MMD (\downarrow) of the adjacency matrices between the generated graphs to the test graphs.

In this experiment, we select the top $[\alpha n]$ eigenvalues to conduct the sampling process with a well-trained GSDM, where α is chosen from $\{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$. The results are shown in Table 3.

Numerical experiments prove that the low-rank α -quantile GSDM retains the state-of-the-art performance among most of the baseline methods. Especially, for the Enzymes dataset, by merely preserving the top-30% eigenvalues and eigenvectors, the 0.3-quantile GSDM can still retain 98.2% of the performance of the full GSDM. The ablation study justifies that our proposed GSDM is capable of capturing essential knowledge that resides on low-dimensional manifolds. It further demonstrates that α -quantile GSDM enables outstanding scalability for large-scale datasets.

Ablation studies on graphs with different eigenvalues distributions. To analyze how GSDM performs under graphs with different eigenvalue distributions, we randomly generate 3 types of synthetic graphs with eigenvalues distributions shown in Figure 6: synthetic-1 has evenly distributed eigenvalues, and synthetic-2 and synthetic-3 have moderate and high distinctions in eigenvalues. Due to the predefined eigenvalue distributions, the edge values of the synthetic graphs are not limited to $\{0, 1\}$. Thus instead of using the metrics presented in Section 5.1, we directly compute the MMD of the adjacency matrices between the generated graphs and the test graphs. Corresponding results compared to GDSS are shown in Figure 7. We also include an ablation variant of GSDM by using only top-50% quantile eigenvalues. Both GSDM and its ablation variant outperform the GDSS model under this synthetic setting.

6 CONCLUSION

In this paper, we present a novel graph diffusion model named GSDM, which consists of diffusion methods on

both graph node features and topologies through SDEs. Specifically, for graph topology diffusion, we design a novel spectral diffusion model in GSDM, to improve the accuracy of predicting score functions, and avoid the pitfalls of the conventional diffusion processes on graphs. Furthermore, we provide theoretical analysis to justify the advantages in effectiveness and efficiency of GSDM compared to the standard graph diffusion. To validate our proposed model, we conducted extensive experiments showing that our proposed GSDM outperforms the state-of-the-art baseline methods on both generic graph generation and molecule generation with significantly higher processing speed.

REFERENCES

- [1] Y. Yang, Z. Feng, M. Song, and X. Wang, "Factorizable graph convolutional networks," in *NeurIPS*, 2020, pp. 20286–20296.
- [2] H. Wang, T. Xu, Q. Liu, D. Lian, E. Chen, D. Du, H. Wu, and W. Su, "Mcne: an end-to-end framework for learning multiple conditional network representations of social network," in *KDD*, 2019, pp. 1064–1072.
- [3] C. Zang and F. Wang, "Moflow: an invertible flow model for generating molecular graphs," in *KDD*, 2020, pp. 617–626.
- [4] C. Shi, M. Xu, Z. Zhu, W. Zhang, M. Zhang, and J. Tang, "Graphaf: a flow-based autoregressive model for molecular graph generation," *arXiv preprint arXiv:2001.09382*, 2020.
- [5] H. Lee, E. Hyung, and S. J. Hwang, "Rapid neural architecture search by learning to generate graphs from datasets," *arXiv preprint arXiv:2107.00860*, 2021.
- [6] H. Liu, L. Jing, J. Wen, P. Xu, J. Wang, J. Yu, and M. K. Ng, "Interpretable deep generative recommendation models," *J. Mach. Learn. Res.*, vol. 22, pp. 202–1, 2021.
- [7] H. Wang, J. Wang, J. Wang, M. Zhao, W. Zhang, F. Zhang, W. Li, X. Xie, and M. Guo, "Learning graph representation with generative adversarial nets," *IEEE Transactions on Knowledge and Data Engineering*, vol. 33, no. 8, pp. 3090–3103, 2019.
- [8] M. Simonovsky and N. Komodakis, "Graphvae: Towards generation of small graphs using variational autoencoders," in *ICANN*, 2018, pp. 412–422.
- [9] J. You, R. Ying, X. Ren, W. Hamilton, and J. Leskovec, "Graphrnn: Generating realistic graphs with deep auto-regressive models," in *ICMLg*, 2018, pp. 5708–5717.
- [10] J. Jo, S. Lee, and S. J. Hwang, "Score-based generative modeling of graphs via the system of stochastic differential equations," *arXiv preprint arXiv:2202.02514*, 2022.
- [11] Y. Song and S. Ermon, *Generative Modeling by Estimating Gradients of the Data Distribution*. Red Hook, NY, USA: Curran Associates Inc., 2019.
- [12] C. Niu, Y. Song, J. Song, S. Zhao, A. Grover, and S. Ermon, "Permutation invariant graph generation via score-based generative modeling," in *AISTATS*, 2020, pp. 4474–4484.
- [13] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," *arXiv preprint arXiv:2011.13456*, 2020.
- [14] X. Guo and L. Zhao, "A systematic survey on deep generative models for graph generation," *arXiv preprint arXiv:2007.06686*, 2020.
- [15] N. De Cao and T. Kipf, "Molgan: An implicit generative model for small molecular graphs," *arXiv preprint arXiv:1805.11973*, 2018.
- [16] T. Ma, J. Chen, and C. Xiao, "Constrained generation of semantically valid graphs via regularizing variational autoencoders," in *NeurIPS*, 2018.
- [17] B. D. Anderson, "Reverse-time diffusion equation models," *Stochastic Processes and their Applications*, vol. 12, no. 3, pp. 313–326, 1982.
- [18] P. Vincent, "A connection between score matching and denoising autoencoders," *Neural Comput.*, vol. 23, no. 7, p. 1661–1674, Jul 2011.
- [19] Y. Li, O. Vinyals, C. Dyer, R. Pascanu, and P. Battaglia, "Learning deep generative models of graphs," *arXiv preprint arXiv:1803.03324*, 2018.
- [20] Y. Luo, K. Yan, and S. Ji, "Graphdf: A discrete flow model for molecular graph generation," in *ICML*, 2021, pp. 7192–7203.

- [21] J. Liu, A. Kumar, J. Ba, J. Kiros, and K. Swersky, "Graph normalizing flows," in *NeurIPS*, 2019.
- [22] M. Liu, K. Yan, B. Oztekin, and S. Ji, "Graphebm: Molecular graph generation with energy-based models," *arXiv preprint arXiv:2102.00546*, 2021.
- [23] R. Ramakrishnan, P. O. Dral, M. Rupp, and O. A. Von Lilienfeld, "Quantum chemistry structures and properties of 134 kilo molecules," *Scientific data*, vol. 1, no. 1, pp. 1–7, 2014.
- [24] J. J. Irwin, T. Sterling, M. M. Mysinger, E. S. Bolstad, and R. G. Coleman, "Zinc: a free tool to discover chemistry for biology," *Journal of chemical information and modeling*, vol. 52, no. 7, pp. 1757–1768, 2012.
- [25] G. Landrum *et al.*, "Rdkit: Open-source cheminformatics software," 2016. [Online]. Available: <http://www.rdkit.org/>
- [26] K. Preuer, P. Renz, T. Unterthiner, S. Hochreiter, and C. Klambauer, "Fréchet chemnet distance: a metric for generative models for molecules in drug discovery," *Journal of chemical information and modeling*, vol. 58, no. 9, pp. 1736–1741, 2018.
- [27] F. Costa and K. De Grave, "Fast neighborhood subgraph pairwise distance kernel," in *ICML*, 2010.
- [28] A. Nichol and P. Dhariwal, "Improved denoising diffusion probabilistic models," *CoRR*, vol. abs/2102.09672, 2021.
- [29] D. P. Kingma, T. Salimans, B. Poole, and J. Ho, "Variational diffusion models," *CoRR*, vol. abs/2107.00630, 2021.
- [30] C. Liu, L. Zhu, and M. Belkin, "Toward a theory of optimization for over-parameterized systems of non-linear equations: the lessons of deep learning," *CoRR*, vol. abs/2003.00307, 2020.

APPENDIX A PROOF DETAILS

A.1 Proof of Proposition 1

Proof. Plugging $\mathbf{A}_t = \mathbf{U}_0 \Lambda_t \mathbf{U}_0^\top$ into (12) yields

$$\begin{aligned} d\mathbf{A}_t &= -\frac{1}{2}\sigma_{\Lambda,t}^2 \mathbf{A}_t dt + \sigma_{\Lambda,t} d\mathbf{M}_t, \\ \mathbf{M}_t &\triangleq \mathbf{U}_0 \mathbf{W}_t^\Lambda \mathbf{U}_0^\top = \mathbf{U}_0 \cdot \text{diag}(\mathbf{B}_t^\Lambda) \cdot \mathbf{U}_0^\top. \end{aligned}$$

Since $(\mathbf{M}_t)_{t \in \mathbb{R}}$ is a linear transformation of $(\mathbf{B}_t^\Lambda)_{t \in \mathbb{R}}$, the standard Brownian motion on \mathbb{R}^n , hence $(\mathbf{M}_t)_{t \in \mathbb{R}}$ is a rank- n centered Gaussian process in $\mathbb{R}^{n \times n}$. Moreover, \mathbf{M}_t is characterized by its covariance kernel $\mathcal{K}(s, t) : [0, 1] \times [0, 1] \mapsto \mathbb{R}^{n \times n \times n \times n}$, which is given by

$$\begin{aligned} \mathcal{K}(s, t)_{i,j,k,l} &\triangleq \text{Cov}(\mathbf{M}_s[i, j], \mathbf{M}_t[k, l]) = \mathbb{E}[\mathbf{M}_s[i, j]\mathbf{M}_t[k, l]] \\ &= \mathbb{E}\left(\sum_{h=1}^n \mathbf{U}_0[i, h]\mathbf{B}_s^\Lambda[h]\mathbf{U}_0[j, h]\right)\left(\sum_{h=1}^n \mathbf{U}_0[k, h]\mathbf{B}_t^\Lambda[h]\mathbf{U}_0[l, h]\right) \\ &= \mathbb{E}\left(\sum_{h=1}^n \mathbf{B}_s^\Lambda[h]\mathbf{B}_t^\Lambda[h]\mathbf{U}_0[i, h]\mathbf{U}_0[j, h]\mathbf{U}_0[k, h]\mathbf{U}_0[l, h]\right) \\ &= \sum_{h=1}^n \mathbb{E}(\mathbf{B}_s^\Lambda[h]\mathbf{B}_t^\Lambda[h])\mathbf{U}_0[i, h]\mathbf{U}_0[j, h]\mathbf{U}_0[k, h]\mathbf{U}_0[l, h] \\ &= \min(s, t) \sum_{h=1}^n \mathbf{U}_0[i, h]\mathbf{U}_0[j, h]\mathbf{U}_0[k, h]\mathbf{U}_0[l, h]. \end{aligned}$$

Notice that (10) is a Ornstein–Uhlenbeck process, which admits a closed-form solution

$$\mathbf{A}_t = \mathbf{A}_0 e^{-\frac{1}{2} \int_0^t \sigma_\tau^2 d\tau} + (1 - e^{-\int_0^t \sigma_\tau^2 d\tau}) \mathbf{W}_1^\Lambda. \quad (19)$$

Hence, plugging (19) into $\mathbf{A}_t = \mathbf{U}_0 \Lambda_t \mathbf{U}_0^\top$ yields

$$\mathbf{A}_t = \mathbf{A}_0 e^{-\frac{1}{2} \int_0^t \sigma_\tau^2 d\tau} + (1 - e^{-\int_0^t \sigma_\tau^2 d\tau}) \mathbf{M}_1. \quad (20)$$

This completes the proof. \square

A.2 Proof of Proposition 2

In preparation of the main proof, we first establish some technical lemmas.

Lemma A.1 (Reconstruction Bound for Generic Diffusion). *We first consider the following oracle reversed time SDE on $(\mathbb{R}^d, \|\cdot\|)$*

$$d\bar{\mathbf{Z}}_t = (\mathbf{f}(\bar{\mathbf{Z}}_t, t) - \sigma_t^2 \nabla_{\mathbf{Z}} \log p_t(\bar{\mathbf{Z}}_t)) d\bar{t} + \sigma_t d\bar{\mathbf{B}}_t, \quad t \in [0, 1],$$

and we define the corresponding estimated reverse time SDE as

$$d\hat{\mathbf{Z}}_t = \left(\mathbf{f}(\hat{\mathbf{Z}}_t, t) - \sigma_t^2 s_\phi(\hat{\mathbf{Z}}_t, t) \right) d\bar{t} + \sigma_t d\bar{\mathbf{B}}_t, \quad t \in [0, t],$$

where $s_\phi(\cdot)$ is optimized to predict the Stein score function $\nabla_{\mathbf{Z}} \log p_t(\mathbf{Z}_t)$ by minimizing the score matching objective

$$\min_{\phi} \mathcal{E}(\phi) \triangleq \mathbb{E}_{\mathbf{Z}} \mathbb{E}_{\mathbf{Z}_t | \mathbf{Z}} \|s_\phi(\mathbf{Z}_t, t) - \nabla_{\mathbf{Z}} \log p_t(\mathbf{Z}_t)\|^2.$$

Then for any ϕ , the construction error is bounded by

$$\begin{aligned} \mathbb{E} \|\mathbf{Z}_0 - \hat{\mathbf{Z}}_0\|^2 &\leq C^2 \|\sigma\|_\infty^4 \mathcal{E}(\phi) \\ &\quad \cdot \left(1 + \int_0^1 F(t) \exp \left(\int_t^1 F(s) ds \right) dt \right), \end{aligned} \quad (21)$$

where $F(t) \triangleq C^2 \sigma_t^4 \|s_\phi(\cdot, t)\|_{\text{lip}}^2 + C \|\mathbf{f}(\cdot, t)\|_{\text{lip}}^2$ and C is a constant.

Proof of Proposition 2. Assumptions of Proposition 2 imply

$$\begin{aligned} \|\mathbf{f}(\cdot, t)\|_{\text{lip}} &= \frac{1}{2} \sigma_t^2, \\ \|s_\phi(\cdot, t)\|_{\text{lip}} &\leq L \cdot \mathbb{E}_{|\mathbf{A}_0} \|\nabla \log p_{t|0}(\hat{\mathbf{A}}_t^{\text{full}})\|_{\text{lip}}, \\ \|s_\phi(\cdot, t)\|_{\text{lip}} &\leq L \cdot \mathbb{E}_{|\mathbf{A}_0} \|\nabla \log p_{t|0}(\hat{\mathbf{A}}_t^{\text{spec}})\|_{\text{lip}}, \end{aligned}$$

where L is a constant. According to [11], [13], at each time step $t \in [0, 1]$, the oracle conditional score function $\nabla \log p_{t|0}(\cdot)$ is equivalent to a mapping that maps the input d -dimensional random variable \mathbf{Z}_t to $\Sigma_t^{-1} \varepsilon$. The shown up ε is a d -dimensional standard Gaussian vector that is independent to \mathbf{Z}_t , and Σ_t is the standard deviation of \mathbf{Z}_t , which is given by $\Sigma_t \triangleq \left(1 - e^{-\int_0^t \sigma_s^2 ds}\right)^{\frac{1}{2}}$. Thus, the expected Lipschitz norm of oracle conditional score function is bounded by

$$\begin{aligned} (\mathbb{E}_{|\mathbf{A}_0} \|\nabla \log p_{t|0}(\mathbf{Z}_t)\|_{\text{lip}})^2 &\leq \mathbb{E}_{|\mathbf{A}_0} \|\nabla \log p_{t|0}(\mathbf{Z}_t)\|_{\text{lip}}^2 \\ &\leq \frac{\Sigma_t^{-2}}{\|\mathbf{Z}_0\|^2} \mathbb{E} \|\varepsilon\|^2 = \frac{\Sigma_t^{-2} d}{\|\mathbf{Z}_0\|^2}. \end{aligned} \quad (22)$$

With Lemma A.1 in hand, we only need to plug (22) into (21), by substituting the notations $(\mathbf{Z}, d, \|\cdot\|)$ with $(\hat{\mathbf{A}}_t^{\text{full}}, n^2, \|\cdot\|_{2,2})$ and $(\hat{\mathbf{A}}_t^{\text{spec}}, n, \|\cdot\|_{2,2})$ to bound the Lipschitz norm of the score networks. This leads us to

$$\begin{aligned} \mathbb{E} \|\mathbf{A}_0 - \hat{\mathbf{A}}_0^{\text{full}}\|^2 &\leq M \mathcal{E}(\phi) \cdot \left(1 + n^2 K \cdot \int_0^1 \Sigma_t^{-2} \exp \left(n^2 K \int_t^1 \Sigma_s^{-2} ds \right) dt \right), \\ \mathbb{E} \|\mathbf{A}_0 - \hat{\mathbf{A}}_0^{\text{spec}}\|^2 &\leq M \mathcal{E}(\phi) \cdot \left(1 + n K \cdot \int_0^1 \Sigma_t^{-2} \exp \left(n K \int_t^1 \Sigma_s^{-2} ds \right) dt \right), \end{aligned}$$

where $M \triangleq C^2 \|\sigma\|_\infty^4$ and $K \triangleq 2ML/\mathbb{E}\|\mathbf{A}_0\|_{2,2}$. For the spectral diffusion part, the final proof step is completed by the fact that

$$\begin{aligned} \mathbb{E}\|\mathbf{A}_0 - \widehat{\mathbf{A}}_0^{\text{spec}}\|^2 &= \mathbb{E}\|\mathbf{U}_0 (\Lambda_0 - \widehat{\Lambda}_0^{\text{spec}}) \mathbf{U}_0^\top\|^2 = \mathbb{E}\|\Lambda_0 - \widehat{\Lambda}_0^{\text{spec}}\|^2, \\ \mathbb{E}\|\mathbf{A}_0\|_{2,2} &= \mathbb{E}\|\mathbf{U}_0 \Lambda_0 \mathbf{U}_0^\top\|_{2,2} = \mathbb{E}\|\Lambda_0\|_{2,2}. \end{aligned}$$

This completes the proof. \square

Proof of Lemma A.1. To bound the expected reconstruction error $\mathbb{E}\|\bar{\mathbf{Z}}_t - \widehat{\mathbf{Z}}_t\|^2$, we first analysis how $\mathbb{E}\|\bar{\mathbf{Z}}_t - \widehat{\mathbf{Z}}_t\|^2$ evolves as time t is reversed from 1 to 0. By definition, it holds that

$$\bar{\mathbf{Z}}_t - \widehat{\mathbf{Z}}_t = \int_1^t \left(\mathbf{f}(\bar{\mathbf{Z}}_s, s) - \mathbf{f}(\widehat{\mathbf{Z}}_s, s) \right) + \sigma_s^2 \left(s_\phi(\widehat{\mathbf{Z}}_s, s) - \nabla_{\mathbf{Z}} \log p_s(\bar{\mathbf{Z}}_s) \right) d\bar{s}.$$

By taking norm and expectation on both sides, we have

$$\begin{aligned} &\mathbb{E}\|\bar{\mathbf{Z}}_t - \widehat{\mathbf{Z}}_t\|^2 \\ &\leq \mathbb{E} \int_1^t \left\| \left(\mathbf{f}(\bar{\mathbf{Z}}_s, s) - \mathbf{f}(\widehat{\mathbf{Z}}_s, s) \right) + \sigma_s^2 \left(s_\phi(\widehat{\mathbf{Z}}_s, s) - \nabla_{\mathbf{Z}} \log p_s(\bar{\mathbf{Z}}_s) \right) \right\|^2 d\bar{s} \\ &\leq C \mathbb{E} \int_1^t \left\| \mathbf{f}(\bar{\mathbf{Z}}_s, s) - \mathbf{f}(\widehat{\mathbf{Z}}_s, s) \right\|^2 d\bar{s} \\ &\quad + C \mathbb{E} \int_1^t \sigma_s^4 \left\| s_\phi(\widehat{\mathbf{Z}}_s, s) - \nabla_{\mathbf{Z}} \log p_s(\bar{\mathbf{Z}}_s) \right\|^2 d\bar{s} \\ &\leq C^2 \int_1^t \sigma_s^4 \cdot \mathbb{E} \left\| s_\phi(\widehat{\mathbf{Z}}_s, s) - s_\phi(\bar{\mathbf{Z}}_s, s) \right\|^2 \\ &\quad + \sigma_s^4 \cdot \mathbb{E} \left\| s_\phi(\bar{\mathbf{Z}}_s, s) - \nabla_{\mathbf{Z}} \log p_s(\bar{\mathbf{Z}}_s) \right\|^2 d\bar{s} \\ &\quad + C \int_1^t \|\mathbf{f}(\cdot, s)\|_{\text{lip}}^2 \cdot \mathbb{E} \left\| \bar{\mathbf{Z}}_s - \widehat{\mathbf{Z}}_s \right\|^2 d\bar{s} \\ &\leq \int_1^t \underbrace{\left(C^2 \sigma_s^4 \|s_\phi(\cdot, s)\|_{\text{lip}}^2 + C \|\mathbf{f}(\cdot, s)\|_{\text{lip}}^2 \right)}_{F(s)} \cdot \mathbb{E} \left\| \widehat{\mathbf{Z}}_s - \bar{\mathbf{Z}}_s \right\|^2 d\bar{s} \\ &\quad + C^2 \underbrace{\int_1^t \sigma_s^4 d\bar{s}}_{G(t)} \cdot \mathcal{E}(\phi). \end{aligned}$$

The proof is completed by applying the Grönwall's inequality to $t \mapsto \mathbb{E}\|\bar{\mathbf{Z}}_t - \widehat{\mathbf{Z}}_t\|^2$, which yields

$$\begin{aligned} &\mathbb{E}\|\mathbf{Z}_0 - \widehat{\mathbf{Z}}_0\|^2 \\ &\leq G(0) + \int_0^1 G(t) F(t) \exp \left(\int_t^1 F(s) ds \right) dt \\ &\leq C^2 \|\sigma\|_\infty^4 \mathcal{E}(\phi) \cdot \left(1 + \int_0^1 F(t) \exp \left(\int_t^1 F(s) ds \right) dt \right). \end{aligned}$$

\square

A.3 Proof of Proposition 3

Proof. Since $\mathcal{E}(\theta) = \mathbb{E}_{S \sim \mathcal{D}_N} \widehat{\mathcal{E}}(\theta; S_N)$, we only need to bound the empirical risk $\widehat{\mathcal{E}}(\theta; S_N)$. By assumption, denote

$$\mathbf{h} \triangleq \begin{pmatrix} s_\theta(\mathbf{Z}_0^1) - \nabla \log p(\mathbf{Z}_0^1 | \mathbf{Z}_1^1)^\top \\ \vdots \\ s_\theta(\mathbf{Z}_0^N) - \nabla \log p(\mathbf{Z}_0^N | \mathbf{Z}_1^N)^\top \end{pmatrix} \in \mathbb{R}^{Nd \times 1},$$

It holds that

$$\|\nabla_\theta \widehat{\mathcal{E}}(\theta; S_N)\|^2 = \frac{1}{N^2} \mathbf{h}^\top \mathbf{K}_\theta(S) \mathbf{h} \geq \frac{\lambda}{N^2} \|\mathbf{h}\|^2 = \frac{\lambda}{N} \widehat{\mathcal{E}}(\theta; S_N),$$

which implies that $\widehat{\mathcal{E}}(\theta; S_N)$ satisfies the $\frac{\lambda}{N}$ -Polyak Łojasiewicz condition [30]. Then the proof is completed by applying Theorem 7 in [30] to $\widehat{\mathcal{E}}(\cdot; S_N)$. \square

APPENDIX B ALGORITHM AND ADDITIONAL RESULTS

Algorithm 1 Training GSDM via minimizing score-matching objective

Input: Score networks $s_{\theta,t}(\cdot)$, $s_{\phi,t}(\cdot)$, maximal diffusion time T , drift functions $\mathbf{f}^X(\cdot, t)$, $\mathbf{f}^\Lambda(\cdot, t)$, noise schedules $\sigma_{X,t}$, $\sigma_{\Lambda,t}$, learning rate η and training epochs K .

Output: Optimized score network parameters θ_K, ϕ_K .

- 1: Initialize θ_0, ϕ_0
- 2: **for** $k = 1$ to K **do**
- 3: $(\mathbf{X}_0, \mathbf{A}_0) \sim \mathcal{G}$
- 4: $\Lambda_0 \leftarrow \text{EigenValues}(\mathbf{A}_0)$
- 5: $t \sim \text{Unif}([0, T])$
- 6: $\mathbf{X}_t \sim \int_0^t f^X(\mathbf{X}_\tau, \tau) d\tau + \int_0^t \sigma_{X,\tau} d\mathbf{B}_\tau^X$, $\Lambda_t \sim \int_0^t f^\Lambda(\Lambda_\tau, \tau) d\tau + \int_0^t \sigma_{\Lambda,\tau} d\mathbf{B}_\tau^\Lambda$
- 7: $\widehat{\mathcal{E}}(\theta_k) \leftarrow \|s_{\theta,k}(\mathbf{X}_t, \Lambda_t) - \nabla \log p_{t|0}(\mathbf{X}_t | \mathbf{X}_0)\|^2$
- 8: $\widehat{\mathcal{E}}(\phi_k) \leftarrow \|s_{\phi,k}(\mathbf{X}_t, \Lambda_t) - \nabla \log p_{t|0}(\Lambda_t | \Lambda_0)\|^2$
- 9: $(\theta_{k+1}, \phi_{k+1}) \leftarrow (\theta_k, \phi_k) - \eta(\nabla \widehat{\mathcal{E}}(\theta_k), \nabla \widehat{\mathcal{E}}(\phi_k))$
- 10: **end for**
- 11: **Return:** θ_K, ϕ_K

Algorithm 2 Sampling via GSDM with VP-SDE predictor-corrector solver

Input: Score-based models $s_{\theta,t}(\cdot)$ and $s_{\phi,t}(\cdot)$, maximal diffusion time T , number of sampling steps M , Langevin-MCMC step sizes $\{\epsilon_i\}_{i=1}^M$, noise schedules $\{\beta_i\}_{i=1}^M$ and priori distribution π .

Output: Generated graph data $(\widehat{\mathbf{X}}_0, \widehat{\mathbf{A}}_0)$

- 1: $t \leftarrow T$
- 2: $(\widehat{\mathbf{X}}_T, \widehat{\mathbf{A}}_T) \sim \pi, \widehat{\mathbf{U}}_0 \sim \text{Unif}(\{\mathbf{U}\} \triangleq \text{EigenVectors}(\mathbf{A}), (\mathbf{X}, \mathbf{A}) \sim \text{Data})$
- 3: **for** $m = M - 1$ to 0 **do**
- 4: $\mathbf{S}_X \leftarrow s_{\theta,t}(\widehat{\mathbf{X}}_t, \widehat{\mathbf{A}}_t, \widehat{\mathbf{U}}_0), \mathbf{S}_\Lambda \leftarrow s_{\phi,t}(\widehat{\mathbf{X}}_t, \widehat{\mathbf{A}}_t, \widehat{\mathbf{U}}_0)$
- 5: $t' \leftarrow t - T/(2M)$
- 6: $\widehat{\mathbf{X}}_{t'} \leftarrow (2 - \sqrt{1 - \beta_{m+1}} \widehat{\mathbf{X}}_t + \beta_{m+1} \mathbf{S}_X) + \sqrt{\beta_{m+1}} \mathbf{z}_X, \mathbf{z}_X \sim N(\mathbf{0}, \mathbf{I})$ {Prediction step: \mathbf{X} }
- 7: $\widehat{\mathbf{A}}_{t'} \leftarrow (2 - \sqrt{1 - \beta_{m+1}} \widehat{\mathbf{A}}_t + \beta_{m+1} \mathbf{S}_\Lambda) + \sqrt{\beta_{m+1}} \mathbf{z}_\Lambda, \mathbf{z}_\Lambda \sim N(\mathbf{0}, \mathbf{I})$ {Prediction step: Λ }
- 8: $\mathbf{S}_X \leftarrow s_{\theta,t'}(\widehat{\mathbf{X}}_{t'}, \widehat{\mathbf{A}}_{t'}, \widehat{\mathbf{U}}_0), \mathbf{S}_\Lambda \leftarrow s_{\phi,t'}(\widehat{\mathbf{X}}_{t'}, \widehat{\mathbf{A}}_{t'}, \widehat{\mathbf{U}}_0)$
- 9: $t \leftarrow t' - T/(2M)$
- 10: $\widehat{\mathbf{X}}_t \leftarrow \widehat{\mathbf{X}}_{t'} + \epsilon_i \mathbf{S}_X + \sqrt{2\epsilon_i} \mathbf{z}_X, \mathbf{z}_X \sim N(\mathbf{0}, \mathbf{I})$ {Correction step: \mathbf{X} }
- 11: $\widehat{\mathbf{A}}_t \leftarrow \widehat{\mathbf{A}}_{t'} + \epsilon_i \mathbf{S}_\Lambda + \sqrt{2\epsilon_i} \mathbf{z}_\Lambda, \mathbf{z}_\Lambda \sim N(\mathbf{0}, \mathbf{I})$ {Correction step: Λ }
- 12: **end for**
- 13: $\widehat{\mathbf{A}}_0 \leftarrow \widehat{\mathbf{U}}_0 \widehat{\mathbf{A}}_0 \widehat{\mathbf{U}}_0^\top$
- 14: **Return:** $(\widehat{\mathbf{X}}_0, \widehat{\mathbf{A}}_0)$

B.1 Algorithm of GSDM

The pseudo codes of training and sampling with GSDM are summarized in Algorithm 1 and Algorithm 2. In a nutshell, we first train a GSDM from real data via minimizing score-matching error. On top of that, we are able to generate graph features and eigen-values of graph adjacency matrices by reversing the forward spectral SDE. By uniformly sampling eigen-vectors from training set, we can construct plausible graph adjacency matrix via spectral composition.

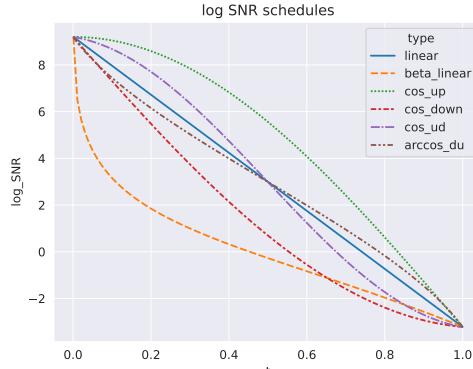


Fig. 8: Diffusion schedules.

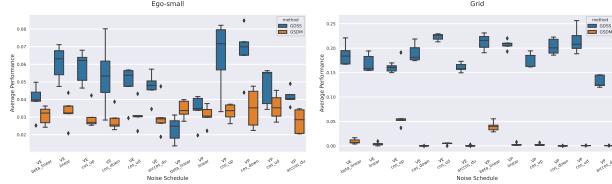


Fig. 9: Ablation study on the choice of noise schedules.

B.2 Ablation study on diffusion schedules

In this section, we empirically verify the robustness of our proposed GSDM under various types of diffusion schedules. As shown in Figure 8, we design six representative noise schedules with the same initial and final signal-to-noise ratio (defined in [13], [28]), while exhibiting different behaviour along the diffusion process. For both VP- and VE-SDE configurations, we train our proposed GSDM with six representative noise schedules and we evaluate the average performance scores. Results on Ego-small and Grid in Figure 9 show that GSDM is robust to the choice of diffusion schedules, and it is able to achieve SOTA performance without deliberate noise schedule optimization.

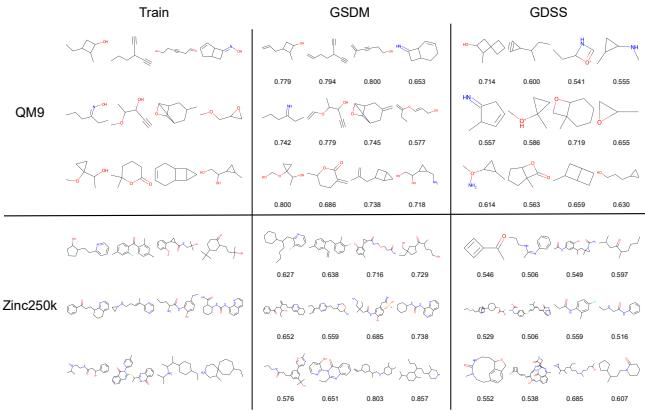


Fig. 10: Visualization of molecule generation with maximum Tanimoto similarity of GSDM comparing to GDSS. The left part shows randomly selected molecules from the training set of QM9 and Zinc250k. For each generated molecules, we show the Tanimoto similarity value at the bottom.

B.3 Visualization of molecule generations

In Figure 10, we show the generated molecules that are maximally similar to certain training molecules. We compute the molecule similarity using the Tanimoto similarity based on the Morgan fingerprints, which are implemented based on REKit [25]. Higher Tanimoto scores indicate that the model is able to generate more similar molecules as the training set, which reflects the learning capability of the model. As shown in the figure, compared to the GDSS model, GSDM is able to generate molecules that have a more similar distribution as the molecules in the training set.

Tianze Luo received the Bachelor degree (First Class Honours) and the Master degree from Nanyang Technological University, Singapore, in 2017 and 2019, respectively. He is currently working towards an Alibaba Talent Programme Ph.D. degree with Alibaba-NTU Joint Research Institute, Nanyang Technological University, Singapore.

PLACE
PHOTO
HERE

Zhanfeng Mo is currently a Ph.D. student in the School of Computer Science and Engineering of Nanyang Technological University, Singapore. He received his Bachelor degree in Statistics from University of Science and Technology of China in 2020. His main research interests include theory and algorithms of statistical machine learning.

PLACE
PHOTO
HERE

Sinno Jialin Pan is a Provost's Chair Professor with Nanyang Technological University (NTU), Singapore. He received his Ph.D. degree in computer science from the Hong Kong University of Science and Technology (HKUST) in 2011. Prior to joining NTU, he was a scientist and Lab Head of text analytics with the Data Analytics Department, Institute for Infocomm Research, Singapore. He serves as an Associate Editor for IEEE TPAMI, AIJ, and ACM TIST.

PLACE
PHOTO
HERE