

BIAS MITIGATION IN GRAPH DIFFUSION MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Most existing graph generative diffusion models suffer from significant exposure bias during graph sampling. We observe that the forward diffusion’s maximum perturbation distribution in most models deviates from the standard normal distribution, while reverse sampling consistently starts from a standard normal distribution. This mismatch results in a reverse starting bias, which, together with the exposure bias, degrades generation quality. The exposure bias typically accumulates and propagates throughout the sampling process. In this paper, we effectively address both biases. To mitigate reverse starting bias, we employ a newly designed Langevin sampling algorithm to align with the forward maximum perturbation distribution, establishing a new reverse starting point. To address the exposure bias, we introduce a fraction correction mechanism based on a newly defined score difference. Our approach, which requires no network modifications, is validated across multiple models, datasets, and tasks, achieving state-of-the-art results.

1 INTRODUCTION

In recent years, graph diffusion models have made significant progress. GDSS (Jo et al., 2022) introduced the score-based diffusion model to the one-shot graph generative task, demonstrating remarkable results and proving superior to baselines. Then, more advanced graph diffusion models such as MOOD (Lee et al., 2023), GSDM (Luo et al., 2023), and HGDM (Wen et al., 2024) were proposed. Given the constraints of graph data scale and network learning capacity, these models truncate the forward diffusion process to enhance performance, preventing it from fully reaching the standard Gaussian distribution. However, during sampling, they have to start from the standard Gaussian distribution without employing any specific strategy. We identify this mismatch as a critical issue and including exposure bias, we work on bias analysis and mitigation in graph diffusion models.

diffusion models (Ho et al., 2020; Song et al., 2021) consist of a forward noising and a reverse denoising process. In the forward process, the data is gradually corrupted by noise over multiple steps. This process can be divided into four stages with the reduction of the signal-to-noise ratio: (1) the data distribution, (2) the low-noise stage, (3) the high-noise stage, and (4) the standard Gaussian.

Reverse-Starting Bias. Ideally, the forward process gradually perturbs the data distribution to the standard Gaussian, while the reverse process starts from the standard Gaussian and gradually recovers clean samples. However, in graph learning, due to limitations in data scale and the network’s learning ability, it is difficult to accurately predict scores from the high-noise state. This forces the forward perturbation to adopt a conservative strategy, where the maximum perturbation distribution falls far short of the standard Gaussian (Jo et al., 2022; Luo et al., 2023; Wen et al., 2024; Lee et al., 2023). Yet, the sampling starting point remains standard Gaussian, resulting in a severe reverse-starting bias, as shown in Figs. 1(a) and 1(b), which significantly affects the generation quality.

Exposure Bias. During the training phase of diffusion model, the model generates corrupted samples \mathbf{x}_t based on ground truth with noise. During the sampling phase, the model starts from a standard Gaussian distribution and iteratively denoises to obtain predicted samples $\hat{\mathbf{x}}_t$ using the score network. Due to the prediction error of the score network, this leads to the exposure bias: a mismatch between \mathbf{x}_t in the training phase and $\hat{\mathbf{x}}_t$ in the sampling phase. This bias accumulates and propagates as sampling progresses, ultimately affecting the generation quality. Naturally, the most direct approach to address exposure bias is to reduce the prediction error of the score network.

Rather than exploring the two biases independently, this paper aims to analyze and mitigate these two biases in graph diffusion models from a unified perspective:

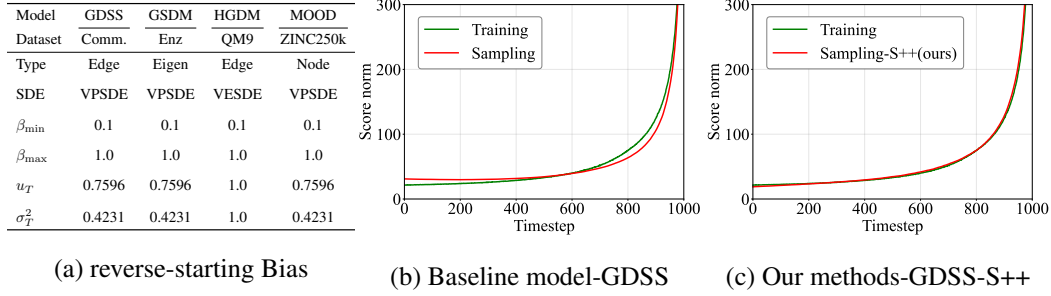


Figure 1: (a) According to the forward formula, we can always write the perturbation distribution as $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t|u_t\mathbf{x}_0, \sigma_t^2\mathbf{I})$. The maximum perturbation distribution of these graph diffusion models in the training phase is $\mathcal{N}(\mathbf{x}_T|u_T\mathbf{x}_0, \sigma_T^2\mathbf{I})$, but the reverse-starting point of the sampling phase always obeys $\mathcal{N}(\mathbf{x}_T|\mathbf{0}, \mathbf{I})$, resulting in significant inconsistencies between training and sampling. In particular, we find that \mathbf{x}_T of different baselines are in the low-noise state, since their signal-to-noise ratio are always greater than 1. We provide more details in appendix A. (b) and (c) Expectation of $\|s_{\theta,t}(\cdot)\|_2$ during sampling (without corrector) and training on Community-small. Due to the initial deviation, there is a significant difference between training and inference in the early stages of sampling in (b). However, after improvement through our method, not only is the reverse-starting bias mitigated, but the exposure bias during the sampling process is also further alleviated.

Q1: Is it possible to mitigate exposure bias while addressing Reverse-Starting bias? It originates from a key finding: when \mathbf{x}_t is in the high-noise stage, the model is highly sensitive to the prediction error of the score network, which means the prediction error at this stage can significantly affect the generative quality. Conversely, when \mathbf{x}_t is in the low-noise state, the model is quite resistant to the prediction error. Coincidentally, the forward maximum perturbation distribution of many models are in the low-noise state as shown in Fig. 1(a). Thus, for a given score network $s_{\theta,t}(\cdot)$, we use Langevin sampling with $s_{\theta,T}(\cdot)$ to obtain samples of the forward maximum perturbation distribution $q(\mathbf{x}_T|\mathbf{x}_0)$. It solves the reverse-starting bias, meanwhile it pushes the reverse-starting point towards the low-noise state, utilizing the model’s resistance to prediction error to avoid exposure bias.

However, the prediction error of the score network severely affects the stable distribution of Langevin sampling, forcing us to improve the prediction accuracy of the network, which is also beneficial for mitigating exposure bias in the sampling process. In particular, we also focus on the cost of achieving:

Q2: How to correct scores without modifying the network or introducing other components?

It originates from a key situation: current graph diffusion models design different networks based on various standards (spatial domain, spectral domain and hyperbolic domain, etc.). We hope our correction method is seamlessly integrated into these models without modifying networks, which means we do not alter the network architecture or model parameters. We also do not introduce any additional components, such as GAN (Goodfellow et al., 2014), Flow (Kingma & Dhariwal, 2018), or discriminator (Kim et al., 2023). We aim to fully utilize existing components of diffusion models to solve their own bias problems. Firstly, we use the pre-trained score network to generate a batch of samples; Then, we train a pseudo score network based on the generated samples; Finally, we use difference of two score networks to correct scores. In summary, our contributions are:

- To the best of our knowledge, we are the first to systematically address bias issues in graph diffusion models, effectively employing Langevin sampling to resolve reverse-starting bias while significantly mitigating exposure bias in the graph sampling.
- We propose a score correction mechanism based on the score difference, and prove both theoretically and practically that the corrected scores are closer to the true scores, further mitigating reverse-starting bias exposure bias.
- Our method does not require modifying the network or introducing new components. It has been validated on multiple graph diffusion models, multiple datasets, and multiple tasks, achieving state-of-the-art metrics.

2 RELATED WORK

Diffusion models were first introduced by Sohl-Dickstein et al. (2015) and later improved by Ho et al. (2020). Notably, Song et al. (2021) proposed a unified framework for diffusion models based on stochastic differential equations, greatly advancing their development. GDSS (Jo et al., 2022) was the first to introduce diffusion models to both nodes and edges of graphs. GSDM (Luo et al., 2023) extended GDSS by introducing the diffusion process of adjacency matrices into the spectral domain. HGDM (Wen et al., 2024) introduced node diffusion into hyperbolic space based on degree distribution characteristics. Huang et al. (2023) proposed a conditional diffusion model based on discrete graph structures. Additionally, Vignac et al. (2023) defined a discrete denoising diffusion model through the process of adding or removing edges and changing categories. Furthermore, Xu et al. (2022) proposed a diffusion model for predicting molecular conformations.

The reverse-starting bias of diffusion models was first discovered by Lin et al. (2024), which proposed modifying the diffusion noise schedule to force the last time step of forward diffusion to have zero Signal-to-Noise Ratio. Shortly after, Everaert et al. (2024) estimated the actual maximum perturbation distribution of forward noise addition as the starting point for inference to match the endpoint of forward noise addition. The exposure bias of diffusion models was first discovered by ADM-IP (Ning et al., 2023), which proposed re-perturbing the perturbation distribution to simulate exposure bias during inference. EB-DDPM (Li & van der Schaar, 2023) estimated the upper bound of cumulative errors and used it as a regularization term to retrain the model. MDSS (Ren et al., 2024) proposed a multi-step timed sampling strategy to mitigate exposure bias. It’s worth noting that ADM-IP, EB-DDPM, and MDSS all require model retraining. In contrast, ADM-ES (Ning et al., 2024) proposed a noise scaling mechanism to mitigate exposure bias without retraining, while TS-DPM (Li et al., 2024) only needs to find the optimal time steps during inference to match the forward process as closely as possible.

We emphasize that our work focuses more on the reverse-starting bias, hoping to address it by utilizing components of diffusion model itself while also mitigating exposure bias to some extent. This is a novel and interesting perspective.

3 MOTIVATION

3.1 GRAPH DIFFUSION MODELS

First, we define a graph with N nodes as $G = (X, A)$, where $X \in \mathbb{R}^{N \times F}$ represents node features, with F indicating that each node has F features; $A \in \mathbb{R}^{N \times N}$ represents the weighted adjacency matrix. Then, we formally represent the graph diffusion process as the trajectory of the random variable G over time $[0, T]$, as shown below:

$$dG_t = f_t(G_t)dt + g_t(G_t)d\mathbf{w}, \quad G_0 \sim p_{\text{data}}. \quad (1)$$

We view this diffusion process as an SDE, where $f_t(G_t)$ is the linear drift coefficient, $g_t(G_t)$ is the diffusion coefficient, \mathbf{w} is a standard Wiener process, and G_0 is a graph from the original distribution p_{data} . Specifically, we replace G in Eq. (1) with node X or edge A , representing the forward diffusion process of node X or edge A , separately.

Following GDSS (Jo et al., 2022), we separate X and A in the reverse diffusion:

$$\begin{aligned} dX_t &= [f_{1,t}(X_t) - g_{1,t}^2 \nabla_{X_t} \log p_t(X_t, A_t)] d\bar{t} + g_{1,t} d\bar{\mathbf{w}}_1, \\ dA_t &= [f_{2,t}(A_t) - g_{2,t}^2 \nabla_{A_t} \log p_t(X_t, A_t)] d\bar{t} + g_{2,t} d\bar{\mathbf{w}}_2 \end{aligned} \quad (2)$$

where $f_{1,t}$ and $f_{2,t}$ satisfy $f_t(X, A) = (f_{1,t}(X), f_{2,t}(A))$, representing the drift coefficients of the reverse-diffusion process for nodes and edges, respectively. $g_{1,t}$ and $g_{2,t}$ are the corresponding scalar diffusion coefficients, $\bar{\mathbf{w}}_1$ and $\bar{\mathbf{w}}_2$ are standard Wiener processes in reverse time, and $\nabla_{X_t} \log p(X_t, A_t)$ and $\nabla_{A_t} \log p(X_t, A_t)$ represent the partial scores of nodes and edges, respectively. It’s worth noting that each SDE in Eq. (2) corresponds to the diffusion process of X and A respectively. We choose different types of SDEs for X and A based on actual conditions. For example, for VPSDE (Song et al., 2021), $f_{1,t}(X_t) = -\frac{1}{2}\beta(t)X_t$, $f_{2,t}(A_t) = -\frac{1}{2}\beta(t)A_t$, $g_{1,t} = g_{2,t} = \sqrt{\beta(t)}$, $\beta(t) = \bar{\beta}_{\min} + t(\bar{\beta}_{\max} - \bar{\beta}_{\min})$, where $\bar{\beta}_{\max}$ and $\bar{\beta}_{\min}$ are hyperparameters.

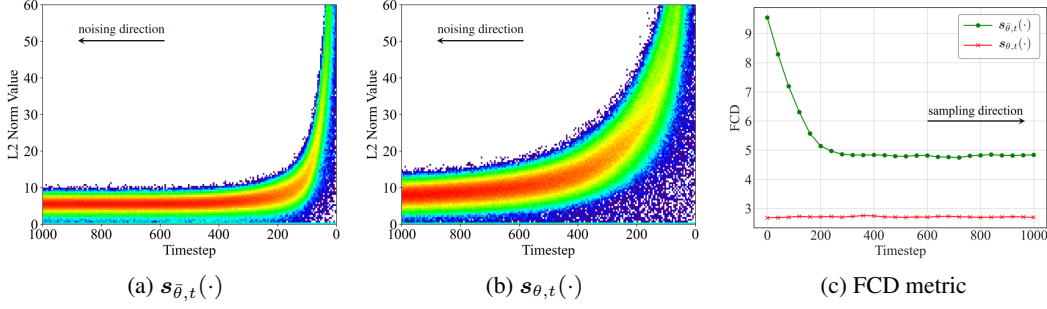


Figure 2: (a) and (b) The ℓ_2 norm distribution of the predicted outputs of the two score networks at different time steps. (c) The response of the predicted outputs of the two score networks to perturbations at different time steps during the sampling phase.

Next, we use $s_{\theta,t}(G_t)$ and $s_{\phi,t}(G_t)$ to estimate the partial scores $\nabla_{X_t} \log p(X_t, A_t)$ and $\nabla_{A_t} \log p(X_t, A_t)$ respectively. Based on the idea of reverse denoising score matching, we derive $s_{\theta,t}(G_t) \approx \nabla_{X_t} \log p_{0t}(X_t|X_0)$ and $s_{\phi,t}(G_t) \approx \nabla_{A_t} \log p_{0t}(A_t|A_0)$, and then give the loss function of the model:

$$\begin{aligned} & \min_{\theta} \mathbb{E}_t \left\{ \lambda_1(t) \mathbb{E}_{G_0} \mathbb{E}_{G_t|G_0} \left\| s_{\theta,t}(G_t) - \nabla_{X_t} \log p_{0t}(X_t|X_0) \right\|_2^2 \right\} \\ & \min_{\phi} \mathbb{E}_t \left\{ \lambda_2(t) \mathbb{E}_{G_0} \mathbb{E}_{G_t|G_0} \left\| s_{\phi,t}(G_t) - \nabla_{A_t} \log p_{0t}(A_t|A_0) \right\|_2^2 \right\} \end{aligned} \quad (3)$$

where $\lambda_1(t)$ and $\lambda_2(t)$ are positive weight functions, t is uniformly sampled from $[0, 1]$. For nodes, we have $X_0 \sim p_0(X)$, $X_t \sim p_{0t}(X_t|X_0)$, and similarly for edges, we have $A_0 \sim p_0(A)$, $A_t \sim p_{0t}(A_t|A_0)$. Since $f_{1,t}$ and $f_{2,t}$ are affine, the transition kernels $p_{0t}(X_t|X_0)$ and $p_{0t}(A_t|A_0)$ are always Gaussian distributions, and closed-form means and variances are obtained based on standard techniques. For example, the node transition kernel in VPSDE (Song et al., 2021) form is shown as follows:

$$p_{0t}(X_t|X_0) = \mathcal{N}\left(X_t \middle| e^{-\frac{1}{4}t^2(\bar{\beta}_{\max}-\bar{\beta}_{\min}) - \frac{1}{2}t\bar{\beta}_{\min}} X_0, \mathbf{I} - \mathbf{I} e^{-\frac{1}{2}t^2(\bar{\beta}_{\max}-\bar{\beta}_{\min}) - t\bar{\beta}_{\min}}\right). \quad (4)$$

For simplicity, the subsequent derivations only focus on X , as the derivations for A are the same as those for X .

3.2 WHY GRAPH DIFFUSION MODELS ARE TRUNCATED?

In this section, we use GDSS as the basic model and QM9 as the dataset to demonstrate the phenomenon of the reverse-starting bias and cleverly corroborate our motivation. We have two score networks: the first is a pretrained network $s_{\theta,t}(\cdot)$ whose forward maximum perturbation is far from reaching standard Gaussian; the second is a network $s_{\bar{\theta},t}(\cdot)$ whose forward maximum perturbation distribution is forced to be standard Gaussian.

Figs. 2(a) and 2(b) show the ℓ_2 norm distribution of the predicted outputs of the two score networks at different time steps. Taking Fig. 2(a) as an example, at each step, we obtain perturbed samples through forward noising, then use $s_{\bar{\theta},t}(\cdot)$ to obtain the predicted score and calculate the corresponding ℓ_2 norm value. We present the details of the figure in Appendix B. At time step 0, the score ℓ_2 norm of the ground truth X_0 spans approximately (0, 2500), demonstrating the diversity of the original data and its scores. As the noise intensity increases, the range of the score ℓ_2 norm narrows, eventually stabilizing within (0, 10). The evolution of the score ℓ_2 norm of perturbed samples at different time steps indicates that as the distribution approaches standard Gaussian, the model becomes highly sensitive to score changes. The tightened score ℓ_2 norm also implies that slight perturbations in scores during the early sampling stages significantly affect generation performance. For Fig. 2(b), the evolution pattern of $s_{\theta,t}(\cdot)$ is consistent with that of $s_{\bar{\theta},t}(\cdot)$, but since the forward perturbation of $s_{\theta,t}(\cdot)$ is far from reaching standard Gaussian, its ℓ_2 norm range is wider, indicating a higher tolerance for score deviations.

Fig. 2(c) illustrates the response of the predicted outputs of the two score networks to perturbations at different time steps during the sampling phase. Each point in Fig. 2(c) represents a

perturbation experiment. The x -axis represents the addition of a standard Gaussian noise perturbation to the score prediction output at the current time step, while the y -axis represents the final generation metric for this perturbation experiment (details in Appendix B). For $s_{\bar{\theta},t}(\cdot)$, at time step 0, we start from standard Gaussian and perturb the predicted score at the current step. Subsequent sampling is not perturbed, ultimately resulting in a rather poor generation metric. As sampling progresses, the destructive effect of score perturbation on generation quality rapidly weakens and stabilizes after 200 steps. The evolution pattern of the score perturbation experiment indicates that diffusion models heavily rely on accurate scores in the early sampling stages, where score deviations severely impact generation quality. For $s_{\theta,t}(\cdot)$, instead of using standard Gaussian as the sampling starting point, we artificially use samples from the forward maximum perturbation distribution as the actual starting point to eliminate the reverse-starting bias.

The above two experiments demonstrate that diffusion models are highly sensitive to score deviations in high-noise states, while in low-noise states, their resistance to score deviations significantly increases. Notably, these experiments also provide us with two directions for addressing the reverse-starting bias: $s_{\bar{\theta},t}(\cdot)$ suggests that we need to retrain and force the forward maximum perturbation distribution to be standard Gaussian, while $s_{\theta,t}(\cdot)$ implies that we need to explore a starting distribution aligned with the forward maximum perturbation distribution during the sampling phase. We find that the latter not only resolves the reverse-starting bias but also provides stronger tolerance to subsequent deviations.

4 METHODOLOGY

4.1 STABLE DISTRIBUTION

Langevin sampling is a key component of SDE-based diffusion models. Given sufficiently small step sizes and a large number of steps, Langevin sampling can utilize the score function to obtain samples from a probability distribution. Importantly, the prior distribution of Langevin sampling can be consistent with that of the diffusion model, typically a standard Gaussian distribution. Moreover, we already have a pretrained score network $s_{\theta,T}(\cdot) \approx \nabla \log q(\mathbf{X}_T|\mathbf{X}_0)$. This score guides Langevin sampling to obtain samples from the distribution $p(\hat{\mathbf{X}}_T) \approx q(\mathbf{X}_T|\mathbf{X}_0)$:

$$\hat{\mathbf{X}}_T^{i+1} \leftarrow \hat{\mathbf{X}}_T^i + \epsilon_T^i s_{\theta}(\hat{\mathbf{X}}_T^i, T) + \sqrt{2\epsilon_T^i} \mathbf{z}_T^i \quad (5)$$

where the subscript T represents the time step parameter of the diffusion model. In the presampling stage, we only use the score $s_{\theta,T}(\cdot)$ at time T . The superscript i denotes the time step parameter of Langevin sampling, ϵ_T^i represents the step size at the current sampling step, and \mathbf{z}_T^i is standard Gaussian noise. After obtaining a batch of samples $\hat{\mathbf{X}}_T$ based on Eq. (5), we use $\hat{\mathbf{X}}_T$ as the new starting point for the reverse sampling process. We refer to this stage as the presampling stage.

4.2 BIAS CORRECTION METHOD

In theory, the presampling stage based on Eq. (5) can obtain samples from the distribution $q(\mathbf{X}_T|\mathbf{X}_0)$. However, the converged score network $s_{\theta,T}(\cdot)$ can never access the true score $\nabla_{\mathbf{X}_T} \log q(\mathbf{X}_T|\mathbf{X}_0) = \frac{\mathbf{X}_T - \sqrt{\bar{\alpha}_T} \mathbf{X}_0}{1 - \bar{\alpha}_T}$. We have to consider the exposure bias of the score network. Without loss of generality,

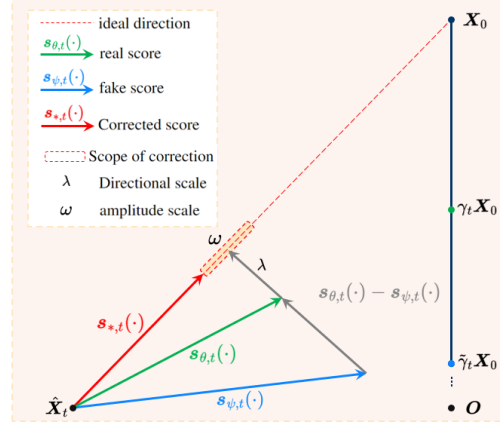


Figure 3: Fractional correction based on the score difference. At the reverse sampling time step t , the ideal score always points to \mathbf{X}_0 . $s_{\theta,t}(\cdot)$ points to $\gamma_t \mathbf{X}_0$ with some deviation (partially containing \mathbf{X}_0), while $s_{\psi,t}(\cdot)$ points to $\tilde{\gamma}_t \mathbf{X}_0$ with larger deviation (containing little \mathbf{X}_0). The difference between real and fake scores guides the real score towards the ideal score. We use λ to control this extent and β to adjust the magnitude of $\hat{s}_{\theta,t}(\cdot)$. The final corrected score flexibly approaches the real score within the dashed box.

we consider the predicted value of the score at any time step:

$$\mathbf{s}_{\theta,t}(\hat{\mathbf{X}}_t) = -\frac{\hat{\mathbf{X}}_t - \sqrt{\bar{\alpha}_t}\hat{\mathbf{X}}_0}{1 - \bar{\alpha}_t}. \quad (6)$$

Clearly, it is challenging for the score network to analytically predict the original data \mathbf{X}_0 . We can rewrite this as:

$$\hat{\mathbf{X}}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(\hat{\mathbf{X}}_t + (1 - \bar{\alpha}_t)\mathbf{s}_{\theta,t}(\hat{\mathbf{X}}_t)). \quad (7)$$

Following Zhang et al. (2023), we model the estimate of $\hat{\mathbf{X}}_0$ as:

$$\hat{\mathbf{X}}_\theta = \gamma_t \mathbf{X}_0 + \eta_t \epsilon_a \quad (8)$$

where $\eta_t < M$, $\epsilon_a \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and for $0 \leq j < k \leq N$, we have $1 > \gamma_j > \gamma_k \geq 0$, $0 \leq \eta_j < \eta_k$. Now, we use the diffusion model to generate a batch of samples $\tilde{\mathbf{X}}_0$ and train a new score network using these as the original data. We know that $\tilde{\mathbf{X}}_0$ always has exposure bias compared to \mathbf{X}_0 , so the score $\mathbf{s}_{\phi,t}(\cdot)$ trained on $\tilde{\mathbf{X}}_0$ naturally learns these exposure bias. Similarly, we define the estimation of the original data by the new network during the reverse sampling, $\hat{\mathbf{X}}_\psi(\hat{\mathbf{X}}_t, t) = \tilde{\gamma}_t \mathbf{X}_0 + \tilde{\eta}_t \epsilon_b$, and we can easily see that $\tilde{\gamma}_t < \gamma_t$. Now we consider the score difference between the two scores at the same time step and for the same sample:

$$\mathbf{s}_{\theta,t}(\hat{\mathbf{X}}_t) - \mathbf{s}_{\psi,t}(\hat{\mathbf{X}}_t) = \sqrt{\bar{\alpha}_t} \frac{(\gamma_t - \tilde{\gamma}_t)\mathbf{X}_0 + (\eta_t \epsilon_a - \tilde{\eta}_t \epsilon_b)}{1 - \bar{\alpha}_t}. \quad (9)$$

We find that the score difference contains information about the original data. We aim to utilize this information. Inspired by classifier-free guidance (Ho & Salimans, 2021) and extrapolation operations (Zhang et al., 2023), we define a new score correction method based on Eq. (9),

$$\begin{aligned} \hat{\mathbf{s}}_{\theta,t}(\hat{\mathbf{X}}_t) &= \mathbf{s}_{\theta,t}(\hat{\mathbf{X}}_t) + \lambda (\mathbf{s}_{\theta,t}(\hat{\mathbf{X}}_t) - \mathbf{s}_{\psi,t}(\hat{\mathbf{X}}_t)) \\ &= -\frac{\hat{\mathbf{X}}_t - \sqrt{\bar{\alpha}_t}((\gamma_t + \lambda(\gamma_t - \tilde{\gamma}_t))\mathbf{X}_0 + \eta_t \epsilon_a + \lambda(\eta_t \epsilon_a - \tilde{\eta}_t \epsilon_b))}{1 - \bar{\alpha}_t} \end{aligned} \quad (10)$$

where $\lambda \geq 0$ represents the step size for correcting the score using the score difference, Eq. (9). When $\lambda = 0$, no correction is applied. Conceptually, the correction operation pulls the biased direction towards the unbiased direction. Although there is some noise in this correction direction, choosing appropriate parameters λ improve the accuracy of the score. Then, we divide the $\hat{\mathbf{s}}_{\theta,t}(\hat{\mathbf{X}}_t)$ by a scalar to adjust the magnitude of the score, further driving $\hat{\mathbf{s}}_{\theta,t}(\hat{\mathbf{X}}_t)$ closer to the true score:

$$\hat{\mathbf{s}}_{\theta,t}(\hat{\mathbf{X}}_t) = \hat{\mathbf{s}}_{\theta,t}(\hat{\mathbf{X}}_t) / \omega. \quad (11)$$

In particular, we emphasize that the score correction at time step T is far more important than at other time steps, as the score is directly related to the steady-state distribution of Langevin sampling, which is crucial for addressing initialization bias. Therefore, we recommend decoupling the correction parameter at time step T from those at other time steps during the actual score correction process.. Specifically, since we use Langevin sampling to obtain a distribution aligned with the forward maximum perturbation distribution, and this distribution is in a low-noise state, retaining some data information from \mathbf{X}_0 , we can shorten the sampling chain, significantly reducing the sampling time. Experimental validation is provided in §5.3.

We emphasize that utilizing Langevin sampling to obtain aligned samples and using the difference signal to correct scores are indispensable components for addressing the reverse-starting bias and the exposure bias. The effect is shown in Fig. 1(c). Additionally, we conduct extensive ablation experiments in §5.4 to demonstrate this point. We provide a detailed geometric illustration in Fig. 3 and provide detailed derivations and proofs of the formulas from §4 in Appendix C.

5 EXPERIMENTS

In this section, we select three generic graph datasets and two molecular datasets to evaluate the performance of our method. In order to demonstrate the broad applicability of this method in addressing

the reverse-starting bias and mitigating exposure bias, we tested it on a variety of mainstream graph diffusion models, namely GDSS (Jo et al., 2022), GSDM (Luo et al., 2023), HGDM (Wen et al., 2024), and MOOD (Lee et al., 2023). Our improved model is prefixed with the basic diffusion model and denoted by S++ at its suffix. At the same time, we perform extensive downstream task testing and ablation study to further illustrate the effectiveness and necessity of S++.

5.1 GENERIC GRAPH GENERATION

Experimental Setup We selected three generic graph datasets to test our approach: (1) Community-small: 100 artificially generated graphs with community structure; (2) Enzymes: 600 protein maps representing the enzyme structure in the BRENDA database (Schomburg et al., 2004); (3) Grid: 100 standard 2D grid diagrams. To evaluate the quality of the generated graphs, we followed the practice of Jo, Lee, and Hwang (2022) and we used the Maximum Mean Difference (MMD) to compare the statistical distribution of the graphs between the same number of generated plots and the test plots, including the distribution of measured degrees, clustering coefficients, and the number of occurrences of the 4-node track.

Dataset Info.	Community-small Synthetic, $12 \leq V \leq 20$				Enzymes Real, $10 \leq V \leq 125$				Grid Synthetic, $100 \leq V \leq 400$			
	Deg.↓	Clus.↓	Orbit↓	Avg.↓	Deg.↓	Clus.↓	Orbit↓	Avg.↓	Deg.↓	Clus.↓	Orbit↓	Avg.↓
GDSS-OC	0.050	0.132	0.011	0.064	0.052	0.627	0.249	0.309	0.270	0.009	0.034	0.070
GDSS-OC-S++	0.021	0.061	0.005	0.029	0.067	0.099	0.007	0.058	0.105	0.004	0.061	0.066
GDSS-WC	0.045	0.088	0.007	0.045	0.044	0.069	0.002	0.038	0.111	0.005	0.070	0.070
GDSS-WC-S++	0.019	0.062	0.004	0.028	0.031	0.050	0.003	0.028	0.105	0.004	0.061	0.057
HGDM-OC	0.065	0.119	0.024	0.069	0.125	0.625	0.371	0.374	0.181	0.019	0.112	0.104
HGDM-OC-S++	0.021	0.034	0.005	0.020	0.080	0.500	0.225	0.268	0.023	0.034	0.004	0.020
HGDM-WC	0.017	0.050	0.005	0.024	0.045	0.049	0.003	0.035	0.137	0.004	0.048	0.069
HGDM-WC-S++	0.021	0.024	0.004	0.016	0.040	0.041	0.005	0.029	0.123	0.003	0.047	0.058
GSDM-OC	0.142	0.230	0.043	0.138	0.930	0.867	0.168	0.655	1.996	0.0	1.013	1.003
GSDM-OC-S++	0.011	0.016	0.001	0.009	0.012	0.087	0.011	0.037	1.2e-4	0.0	1.2e-4	0.066
GSDM-WC	0.011	0.016	0.001	0.009	0.013	0.088	0.013	0.038	0.002	0.0	0	7.2e-5
GSDM-WC-S++	0.011	0.016	0.001	0.009	0.011	0.086	0.010	0.036	5.0e-5	0.0	1.1e-5	0.066

Table 1: Generation results on the generic graph datasets (Lower is better). The results of the Enzymes dataset of GDSS are reproduced by ourselves, and the results of other baselines are all from published papers, and we give detailed settings and instructions in the appendix D.

Results Table 1 shows that S++ significantly outperforms all baseline models. For the uncorrected sampling method, the performance indicators of the baseline model are particularly poor due to the existence of the reverse-starting bias and score exposure bias, while S++ can significantly improve the performance of all baseline models and reach or even exceed the level of the baseline model with correctors. Because the method without correctors can significantly reduce the computational consumption, we believe that S++ can really release the ability of the graph diffusion model, which is enlightening for large-scale datasets. For the sampling method with aligners, S++ is still significantly better than all baseline models. At the same time, we also give experimental comparisons of other advanced models in the appendix F, and the results show that S++ can achieve the SOTA indicators of the corresponding tasks.

5.2 MOLECULAR GRAPH GENERATION

Experimental Setup We selected two widely recognized molecular datasets to evaluate our methods: QM9 (Ramakrishnan et al., 2014) and ZINC250k (Irwin et al., 2012). We generated 10,000 molecules and selected the following widely used evaluation metrics: Frechet ChemNet Distance (FCD) (Preuer et al., 2018), Neighborhood subgraph pairwise distance kernel (NSPDK) MMD (Costa & Grave, 2010), validity w/o correction, and the generation time. (1) FCD uses the activation of the penultimate layer of ChemNet to calculate the distance between the benchmark molecular dataset and the generated dataset to characterize the similarity between the two, and the lower the FCD value, the higher the similarity between the two distributions. (2) (NSPDK) MMD considered the characteristics of nodes and edges at the same time, and calculated the MMD between the benchmark molecular dataset and the generated dataset; (3) Sampling time is used to evaluate the rapidity of the model in generating

large-scale molecular datasets, and we only count the time spent on sampling, regardless of the time spent on preprocessing and evaluation.

Results Table 2 shows that both in terms of sampling time and generation quality, S++ is significantly better than the baseline model. For the sampling method without correctors, due to the existence of the reverse-starting bias and score exposure bias, the quality of generation from the baseline model is particularly poor, while S++ can significantly improve the performance of all baselines and approximate the sampling methods with correctors of the baseline model. For the sampling method with aligners, S++ is still significantly better than all baseline models and greatly reduces the sampling time. At the same time, we provide more comparative experimental results in appendix F and provide parameter sensitivity experiments in appendix G.

Method	QM9			ZINC250k		
	Sampling time ↓	NSPDK MMD ↓	FCD ↓	Sampling time s↓	NSPDK MMD ↓	FCD ↓
GDSS-OC	$0.73e^2$	0.016	4.584	$0.73e^3$	0.047	20.53
GDSS-OC-S++	5.10	0.001	1.661	0.70e³	0.050	16.79
GDSS-WC	$1.61e^2$	0.004	2.550	$1.41e^3$	0.019	14.66
GDSS-WC-S++	9.25	0.001	1.661	0.98e³	0.012	12.70
HGDM-OC	$0.62e^2$	0.005	3.164	$0.76e^3$	0.033	21.38
HGDM-OC-S++	0.62e²	0.003	2.512	0.77e³	0.034	20.79
HGDM-WC	$1.16e^2$	0.002	2.147	$1.52e^3$	0.016	17.69
HGDM-WC-S++	0.98e²	0.001	2.001	1.17e³	0.016	16.24

Table 2: Comparison of different methods on QM9 and ZINC250k datasets.

5.3 DIVERSITY GENERATION

Characteristic molecule generation To evaluate the performance of S++ in generating novel, drug-like, and synthesizable molecules, we follow (Lee et al., 2023) and assess S++ in the five docking score (DS) optimization tasks under the quantitative estimate of synthetic accessibility (SA), drug-likeness (QED) and novelty constraints. We define the property Y by

$$Y(G) = \widehat{DS}(G) \times \widehat{QED}(G) \times \widehat{SA}(G) \in [0, 1] \quad (12)$$

where \widehat{DS} refers to the normalized docking score, \widehat{SA} denotes the normalized synthetic accessibility, and QED represents drug-likeness. We used MOOD-S++ to generate 3000 molecules and evaluate performance using the following metrics. **Novel hit ratio (%)** is the fraction of unique hit molecules whose maximum Tanimoto similarity with the training molecules is less than 0.4. In particular, hit molecules are defined as the molecules that satisfy the following conditions: $DS < (\text{the median DS of the known active molecules})$, $QED > 0.5$, and $SA < 5$. **Novel top 5% docking score** refers to the average DS of the top 5% unique molecules that satisfy the constraints $QED > 0.5$ and $SA < 5$ and their maximum similarity with the training molecules is below 0.4. To avoid bias in target selection, we utilize five protein targets: parp1, fa7, 5ht1b, braf, and jak2.

Results Tables 3 and 4 show that MOOD-S++ is significantly better than baseline in all target proteins. This indicates that S++ still has advantages in the discovery of drug-like, synthesizable, and novel molecular tasks with high binding affinity, and it can be seen that the reverse-starting bias and exposure bias pose a significant threat to various generation tasks.

Accelerate generation To demonstrate that S++ can generate good samples faster by using fewer steps of reverse diffusion, We chose GDSS-OC as the benchmark model, and QM9 and Comm datasets were selected to test the performance of our method and benchmark model at different sampling total time steps.

Results Table 5 shows that S++ is significantly better than the baseline model at different sampling total time steps. S++ was not only able to generate samples with fewer reverse-diffusion steps, but also achieved consistent improvements across generation metrics, especially on the QM9 dataset, where S++ remained close to optimal performance even with a significant reduction in the sampling time step ($T = 100$), while the performance of the benchmark model decreased significantly.

In conclusion, S++ shows higher efficiency, better quality, and stronger robustness in graph generative tasks, which provides a powerful improvement scheme for the application of diffusion model.

Method	Target protein				
	parp1	fa7	5ht1b	braf	jak2
MOOD	7.017 (± 0.428)	0.733 (± 0.141)	18.673 (± 0.423)	5.240 (± 0.285)	9.200 (± 0.524)
MOOD-S++	8.286 (± 0.214)	0.900 (± 0.068)	20.354 (± 0.672)	5.653 (± 0.073)	9.167 (± 0.067)

Table 3: Novel hit ratio (%) results (\uparrow).

Method	Target protein				
	parp1	fa7	5ht1b	braf	jak2
MOOD	-10.865 (± 0.113)	-8.160 (± 0.071)	-11.145 (± 0.042)	-11.063 (± 0.034)	-10.147 (± 0.060)
MOOD-S++	-10.961 (± 0.027)	-8.182 (± 0.028)	-11.231 (± 0.036)	-11.143 (± 0.025)	-10.163 (± 0.015)

Table 4: Novel top 5% docking score (kcal/mol) results (\downarrow).

T	Method	QM9			Community-small			
		Val. w/o corr. \uparrow	NSPDK MMD \downarrow	FCD \downarrow	Deg. \downarrow	Clus. \downarrow	Orbit \downarrow	Avg. \downarrow
1000	GDSS-OC	73.5	0.015	4.584	0.050	0.132	0.011	0.064
	GDSS-OCS++	94.0	0.001	1.671	0.021	0.061	0.005	0.029
500	GDSS-OC	46.2	0.045	7.960	0.136	0.456	0.151	0.248
	GDSS-OC-S++	93.9	0.001	1.665	0.029	0.142	0.008	0.060
100	GDSS-OC	37.8	0.069	9.951	0.092	0.666	0.394	0.384
	GDSS-OC-S++	93.9	0.001	1.663	0.061	0.414	0.140	0.205

Table 5: Comparison of different methods on QM9 and ZINC250k datasets under different total sampling time steps.

Method	QM9		
	Val. w/o corr. \uparrow	NSPDK MMD \downarrow	FCD \downarrow
GDSS-OC	73.5	0.0157	4.58
GDSS-w/o Score Correction	94.8	0.0037	2.65
GDSS-w/o Langevin Alignment	89.8	0.0031	2.01
GDSS-OC-S++	94.0	0.0014	1.67

Table 6: Ablation experiments on the OM9 dataset.

5.4 ABLATION STUDY

Table 6 clearly shows that GDSS-w/o Langevin Alignment or Langevin Alignment alone can improve the performance of the baseline model to varying degrees, however, when we combine these two methods, the model performance is significantly improved, which strongly proves the effectiveness and necessity of the combination of the two methods, and their synergistic effect. Moreover, we provide a comparative analysis of the two biases in the image and graph fields in appendix H, and provide comparative experiments of S++ with existing methods on images in appendix I.

6 CONCLUSION

In this paper, we use Langevin sampling to obtain samples aligned with the forward maximum perturbation distribution, which solves the reverse-starting bias and greatly alleviates the exposure bias of the score network, and we propose a score correction mechanism based on score difference to further promote the stable-state distribution of Langevin sampling to the real forward maximum perturbation distribution, and further alleviate the exposure bias during the sampling phase. Our approach does not require network modifications or the introduction of new components, and can be naturally integrated into existing graph diffusion models to achieve state-of-the-art metrics on multiple datasets and multiple tasks.

REFERENCES

- Fabrizio Costa and Kurt De Grave. Fast neighborhood subgraph pairwise distance kernel. In *ICML*, pp. 255–262, 2010.
- Martin Nicolas Everaert, Athanasios Fitsios, Marco Bocchio, Sami Arpa, Sabine Süsstrunk, and Radhakrishna Achanta. Exploiting the signal-leak bias in diffusion models. In *WACV*, pp. 4025–4034, 2024.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.
- Han Huang, Leilei Sun, Bowen Du, and Weifeng Lv. Conditional diffusion based on discrete graph structures for molecular graph generation. In *AAAI*, 2023.
- John J Irwin, Teague Sterling, Michael M Mysinger, Erin S Bolstad, and Ryan G Coleman. Zinc: a free tool to discover chemistry for biology. *Journal of chemical information and modeling*, pp. 1757–1768, 2012.
- Jaehyeong Jo, Seul Lee, and Sung Ju Hwang. Score-based generative modeling of graphs via the system of stochastic differential equations. In *ICML*, pp. 10362–10383, 2022.
- Dongjun Kim, Yeongmin Kim, Se Jung Kwon, Wanmo Kang, and Il-Chul Moon. Refining generative process with discriminator guidance in score-based diffusion models. In *ICML*, 2023.
- Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *NeurIPS*, 2018.
- Seul Lee, Jaehyeong Jo, and Sung Ju Hwang. Exploring chemical space with score-based out-of-distribution generation. In *ICML*, 2023.
- Mingxiao Li, Tingyu Qu, Ruicong Yao, Wei Sun, and Marie-Francine Moens. Alleviating exposure bias in diffusion models through sampling with shifted time steps. In *ICLR*, 2024.
- Yangming Li and Mihaela van der Schaar. On error propagation of diffusion models. In *ICLR*, 2023.
- Yujia Li, Oriol Vinyals, Chris Dyer, Razvan Pascanu, and Peter Battaglia. Learning deep generative models of graphs. *arXiv preprint arXiv:1803.03324*, 2018.
- Shanchuan Lin, Bingchen Liu, Jiashi Li, and Xiao Yang. Common diffusion noise schedules and sample steps are flawed. In *WACV*, 2024.
- Meng Liu, Keqiang Yan, Bora Oztekin, and Shuiwang Ji. Graphebm: Molecular graph generation with energy-based models. In *Energy Based Models Workshop-ICLR 2021*, 2021.
- Tianze Luo, Zhanfeng Mo, and Sinno Jialin Pan. Fast graph generation via spectral diffusion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- Mang Ning, Enver Sangineto, Angelo Porrello, Simone Calderara, and Rita Cucchiara. Input perturbation reduces exposure bias in diffusion models. In *ICML*, 2023.
- Mang Ning, Mingxiao Li, Jianlin Su, Albert Ali Salah, and Itir Onal Ertugrul. Elucidating the exposure bias in diffusion models. In *ICLR*, 2024.
- Chenhao Niu, Yang Song, Jiaming Song, Shengjia Zhao, Aditya Grover, and Stefano Ermon. Permutation invariant graph generation via score-based generative modeling. In *AISTATS*, 2020.
- Kristina Preuer, Philipp Renz, Thomas Unterthiner, Sepp Hochreiter, and Gunter Klambauer. Fréchet chemnet distance: a metric for generative models for molecules in drug discovery. *Journal of chemical information and modeling*, pp. 1736–1741, 2018.

- Raghunathan Ramakrishnan, Pavlo O Dral, Matthias Rupp, and O Anatole Von Lilienfeld. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific data*, 2014.
- Zhiyao Ren, Yibing Zhan, Liang Ding, Gaoang Wang, Chaoyue Wang, Zhongyi Fan, and Dacheng Tao. Multi-step denoising scheduled sampling: Towards alleviating exposure bias for diffusion models. In *AAAI*, pp. 4667–4675, 2024.
- Ida Schomburg, Antje Chang, Christian Ebeling, Marion Gremse, Christian Heldt, Gregor Huhn, and Dietmar Schomburg. Brenda, the enzyme database: updates and major new developments. *Nucleic acids research*, pp. D431–D433, 2004.
- Chence Shi, Minkai Xu, Zhaocheng Zhu, Weinan Zhang, Ming Zhang, and Jian Tang. GraphAF: a flow-based autoregressive model for molecular graph generation. In *ICLR*, 2020.
- Martin Simonovsky and Nikos Komodakis. Graphvae: Towards generation of small graphs using variational autoencoders. In *ICANN*, pp. 412–422, 2018.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, pp. 2256–2265, 2015.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021.
- Clement Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pascal Frossard. Digress: Discrete denoising diffusion for graph generation. In *ICLR*, 2023.
- Lingfeng Wen, Xuan Tang, Mingjie Ouyang, Xiangxiang Shen, Jian Yang, Daxin Zhu, Mingsong Chen, and Xian Wei. Hyperbolic graph diffusion model. In *AAAI*, pp. 15823–15831, 2024.
- Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. Geodiff: a geometric diffusion model for molecular conformation generation. In *ICLR*, 2022.
- Jiaxuan You, Rex Ying, Xiang Ren, William Hamilton, and Jure Leskovec. Graphrnn: Generating realistic graphs with deep auto-regressive models. In *International conference on machine learning*, 2018.
- Guoqiang Zhang, Kenta Niwa, and W Bastiaan Kleijn. Lookahead diffusion probabilistic models for refining mean estimation. In *CVPR*, 2023.

A REVERSE-STARTING BIAS

In this section, we provide a detailed discussion of the initialization bias in diffusion models. It is worth noting that these diffusion models are based on diffusion models defined by SDE (Song et al., 2021). For VPSDE, the diffusion model obtains perturbed samples through $p_{0i}(\mathbf{X}_i|\mathbf{X}_0) = \mathcal{N}(\mathbf{X}_i|\sqrt{\alpha_i}\mathbf{X}_0, (1 - \alpha_i)\mathbf{I})$, where $\alpha_i := \prod_{j=1}^i(1 - \beta_j)$. When this expression is extended continuously, it leads to Eq. (4), which corresponds to Equation (33) in the SDE. At $t = 1$, Eq. (4) gives the maximum perturbation distribution, which is $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Similarly, for VESDE, the diffusion model obtains perturbed samples through $p_{0i}(\mathbf{X}_i|\mathbf{X}_0) = \mathcal{N}(\mathbf{X}_i|\mathbf{X}_0, \sigma_i\mathbf{I})$, where $\sigma_{\min} = \sigma_1 < \sigma_2 < \dots < \sigma_N = \sigma_{\max}$. When this expression is extended continuously,

$$p_{0t}(\mathbf{X}_t|\mathbf{X}_0) = \mathcal{N}\left(\mathbf{X}_t|\mathbf{X}_0, \sigma_{\min}^2\left(\frac{\sigma_{\max}}{\sigma_{\min}}\right)^{2t}\mathbf{I}\right), \quad (13)$$

it corresponds to Eq. (31) of the paper of SDE (Song et al., 2021). At $t = 1$ Eq. (13) achieves the maximum perturbation distribution, which is $\mathcal{N}(\mathbf{X}_i|\mathbf{X}_0, \sigma_{\max}\mathbf{I})$. In particular, we need to make sure that σ_{\max} is large enough that $\mathcal{N}(\mathbf{X}_i|\mathbf{X}_0, \sigma_{\max}\mathbf{I}) \approx \mathcal{N}(\mathbf{X}_i|\mathbf{0}, \sigma_{\max}\mathbf{I})$.

However, in practice, Lots of diffusion models (Jo et al., 2022; Luo et al., 2023; Wen et al., 2024) adopted a rather conservative strategy when training the network. For VPSDE, this results in the maximum forward perturbation distribution being $\mathcal{N}(\mathbf{X}_T|u_T\mathbf{x}_0, \sigma_T\mathbf{I})$, which is far from reaching $\mathcal{N}(\mathbf{0}, \mathbf{I})$. For VESDE, due to σ_{\max} not being large enough, the maximum forward perturbation distribution is $\mathcal{N}(\mathbf{X}_i|\mathbf{X}_0, \sigma_{\max}\mathbf{I})$, which cannot be approximated by $\mathcal{N}(\mathbf{X}_i|\mathbf{0}, \sigma_{\max}\mathbf{I})$. However, GDSS et al. always start reverse sampling from the standard Gaussian distribution, which leads to significant initialization bias. A detailed comparison of the parameters is shown in Tables 7, 8, and 9.

Model	GDSS									
Dataset	Community-small		Enzymes		Grid		QM9		ZINC250k	
Type	Node	Edge	Node	Edge	Node	Edge	Node	Edge	Node	Edge
SDE	VPSDE	VPSDE	VPSDE	VESDE	VPSDE	VPSDE	VESDE	VESDE	VPSDE	VESDE
β_{\min}	0.1	0.1	0.1	0.2	0.1	0.2	0.1	0.1	0.1	0.2
β_{\max}	1.0	1.0	1.0	1.0	1.0	0.8	1.0	1.0	1.0	1.0
u_T	0.7596	0.7596	0.7596	1.0	0.7596	0.7788	1.0	1.0	0.7596	1.0
σ_T^2	0.4231	0.4231	0.4231	1.0	0.4231	0.3935	1.0	1.0	0.4231	1.0

Table 7: The actual parameters of the forward perturbation of the GDSS.

Model	HGDM									
Dataset	Community-small		Enzymes		Grid		QM9		ZINC250k	
Type	Node	Edge	Node	Edge	Node	Edge	Node	Edge	Node	Edge
SDE	VPSDE	VPSDE	VPSDE	VESDE	VPSDE	VESDE	VPSDE	VESDE	VPSDE	VESDE
β_{\min}	0.1	0.1	0.1	0.2	0.1	0.2	0.1	0.1	0.1	0.2
β_{\max}	1.0	1.0	1.0	1.0	7.0	0.8	2.0	1.0	1.0	1.0
u_T	0.7596	0.7596	0.7596	1.0	0.1695	1.0	0.5916	1.0	0.7596	1.0
σ_T^2	0.4231	0.4231	0.4231	1.0	0.9713	0.64	0.6501	1.0	0.4231	1.0

Table 8: The actual parameters of the forward perturbation of the HGDM.

Model	GSDM						MOOD			
Dataset	Community-small		Enzymes		Grid		QM9		ZINC250k	
Type	Node	Eigen	Node	Eigen	Node	Eigen	Node	Edge	Node	Edge
SDE	VPSDE	VPSDE	VPSDE	VPSDE	VPSDE	VPSDE	VPSDE	VESDE	VPSDE	VESDE
β_{\min}	0.1	0.1	0.1	0.1	0.1	0.2	0.1	0.1	0.1	0.2
β_{\max}	1.0	1.0	1.0	1.0	1.0	0.8	1.0	1.0	1.0	1.0
u_T	0.7596	0.7596	0.7596	0.7596	0.7596	0.7788	1.0	1.0	0.7596	1.0
σ_T^2	0.4231	0.4231	0.4231	0.4231	0.4231	0.3935	1.0	1.0	0.4231	1.0

Table 9: The actual parameters of the forward perturbation of the GSDM and MOOD.

B FIGURE DETAILS

In this section, we present the detailed procedures to plot Fig. 2. Let $s_{\theta,t}(\cdot)$ represent the GDSS pretrained score network. Due to the conservative strategy of GDSS, with $\beta_{\min} = 0.1$ and $\beta_{\max} = 1$, the maximum perturbation distribution is $\mathcal{N}(\mathbf{X}_T | 0.7596\mathbf{X}_0, 0.4231\mathbf{I})$ at $t = 1$. On the other hand, $s_{\psi,t}(\cdot)$ is defined with the forced constraints of $\beta_{\min} = 0.1$ and $\beta_{\max} = 20$. At $t = 1$, the maximum perturbation distribution is $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Then, $s_{\theta,t}(\cdot)$ and $s_{\psi,t}(\cdot)$ not only represent two different networks, but also indicate that their maximum perturbation distributions during the training phases are completely different.

To plot Figs. 2a and 2b, we freeze the converged $s_{\theta,t}(\cdot)$ and $s_{\psi,t}(\cdot)$, then replace \mathbf{X} in Eq. (3) with \mathbf{A} to obtain perturbation samples of 1024 edges at different timesteps. We then compute $\|s_{\theta,t}(\cdot)\|_2$ and $\|s_{\psi,t}(\cdot)\|_2$ and plot them on the figure.

To plot Fig. 2c, we introduce perturbations to $s_{\theta,t}(\cdot)$ at different timesteps during the sampling phase. We employ a sampling method without a corrector and perturb the score at the selected timestep (horizontal axis) using Gaussian noise:

$$s_{\theta,t}(\cdot) = s_{\theta,t}(\cdot) + z_t \quad (14)$$

where $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. For the other timesteps, we do not introduce any perturbations, allowing the diffusion model to perform sampling and record the generation metrics. We conduct the perturbation experiment on $s_{\theta,t}(\cdot)$ using the same method, and ultimately compare the results of the two perturbation experiments based on the timesteps to evaluate how different score networks in the diffusion model resist bias at various timesteps. We present a detailed comparison of the generation metrics from the perturbation experiments, as shown in Fig. 4.

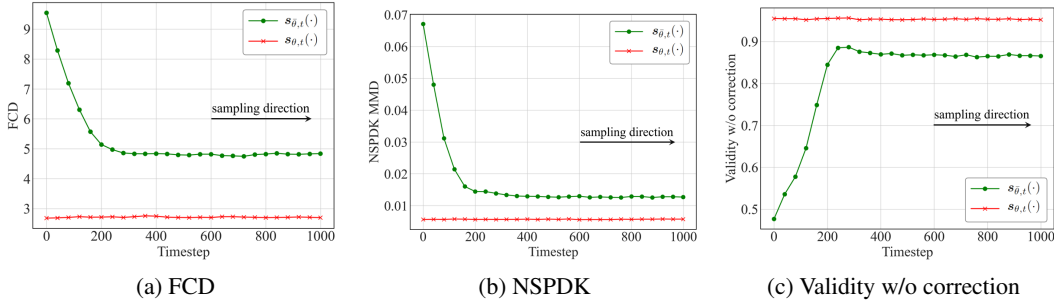


Figure 4: Generation metric responses to perturbations at different timesteps for two-score networks.

C DERIVATIONS FOR §4.2

For a diffusion model, let the original data be \mathbf{X}_0 , and the pretrained score network be $s_{\theta,t}(\cdot)$. Based on the set noise addition method, we have:

$$\nabla_{\mathbf{X}_t} \log q(\mathbf{X}_t | \mathbf{X}_0) = -\frac{\mathbf{X}_t - \sqrt{\bar{\alpha}_t} \mathbf{X}_0}{1 - \bar{\alpha}_t} \quad (15)$$

However, $s_{\theta,t}(\cdot)$ often deviates from the ideal logarithmic gradient. In the reverse process, assuming the current time step t has a data state $\hat{\mathbf{X}}_t$, the score network’s predicted output is:

$$s_{\theta,t}(\hat{\mathbf{X}}_t) = -\frac{\hat{\mathbf{X}}_t - \sqrt{\bar{\alpha}_t} \mathbf{X}_0}{1 - \bar{\alpha}_t} \quad (16)$$

Since it’s difficult for the trained score to analytically predict \mathbf{X}_0 . We model $\hat{\mathbf{X}}_0$ by:

$$\hat{\mathbf{X}}_0 = \gamma_t \mathbf{X}_0 + \eta_t \epsilon_a \quad (17)$$

Eq. (16) becomes:

$$s_{\theta,t}(\hat{\mathbf{X}}_t) = -\frac{\hat{\mathbf{X}}_t - \sqrt{\bar{\alpha}_t}(\gamma_t \mathbf{X}_0 + \eta_t \epsilon_a)}{1 - \bar{\alpha}_t} \quad (18)$$

Now, we train a new score network $s_{\psi,t}(\cdot)$ based on the generated data \bar{X} from the score network. Following the above derivation, we can write the predicted score at the current time step t as:

$$s_{\psi,t}(\hat{X}_t) = -\frac{\hat{X}_t - \sqrt{\bar{\alpha}_t}\tilde{X}_0}{1 - \bar{\alpha}_t} \quad (19)$$

Due to the bias of $s_{\theta,t}(\cdot)$, the generated data X_1 always deviates from X_0 , and considering the prediction error of the network, we can easily obtain:

$$\tilde{X}_0 = \tilde{\gamma}_t X_0 + \tilde{\eta}_t \epsilon_b \quad (20)$$

where $\tilde{\gamma}_t < \gamma_t$, $\tilde{\eta}_t > \eta_t$, meaning that at the same reverse time step t , $s_{\psi,t}(\cdot)$ trained on generated data has a larger bias in predicting the target distribution than $s_{\theta,t}(\cdot)$ trained on original data. We can rewrite $s_{\psi,t}(\cdot)$ as:

$$s_{\psi,t}(\hat{X}_t) = -\frac{\hat{X}_t - \sqrt{\bar{\alpha}_t}(\tilde{\gamma}_t X_0 + \tilde{\eta}_t \epsilon_b)}{1 - \bar{\alpha}_t} \quad (21)$$

Next, we derive the meaning of the score difference, defined as Eq. (18) minus Eq. (21):

$$\begin{aligned} s_{\theta,t}(\hat{X}_t) - s_{\psi,t}(\hat{X}_t) &= -\frac{\hat{X}_t - \sqrt{\bar{\alpha}_t}(\gamma_t X_0 + \eta_t \epsilon_a)}{1 - \bar{\alpha}_t} - \left(-\frac{\hat{X}_t - \sqrt{\bar{\alpha}_t}(\tilde{\gamma}_t X_0 + \tilde{\eta}_t \epsilon_b)}{1 - \bar{\alpha}_t} \right) \\ &= \sqrt{\bar{\alpha}_t} \frac{(\gamma_t - \tilde{\gamma}_t)X_0 + (\eta_t \epsilon_a - \tilde{\eta}_t \epsilon_b)}{1 - \bar{\alpha}_t} \end{aligned} \quad (22)$$

We add this score difference as a correction term to the original predicted score and introduce a hyperparameter to control the influence of the original and the noise scores:

$$\begin{aligned} \hat{s}_{\theta,t}(\hat{X}_t) &= s_{\theta,t}(\hat{X}_t) + \lambda (s_{\theta,t}(\hat{X}_t) - s_{\psi,t}(\hat{X}_t)) \\ &= -\frac{\hat{X}_t - \sqrt{\bar{\alpha}_t}(\gamma_t X_0 + \eta_t \epsilon_a + \lambda(\gamma_t - \tilde{\gamma}_t)X_0 + \lambda(\eta_t \epsilon_a - \tilde{\eta}_t \epsilon_b))}{1 - \bar{\alpha}_t} \\ &= -\frac{\hat{X}_t - \sqrt{\bar{\alpha}_t}((\gamma_t + \lambda(\gamma_t - \tilde{\gamma}_t))X_0 + \eta_t \epsilon_a + \lambda(\eta_t \epsilon_a - \tilde{\eta}_t \epsilon_b))}{1 - \bar{\alpha}_t} \end{aligned} \quad (23)$$

Because $\tilde{\gamma}_t < \gamma_t$, this score difference helps the original score add more information from the original data X_0 . By setting an appropriate hyperparameter λ , we can always use the information from the original score to guide the correction of the score. Finally, we add a coefficient of adjustment amplitude to the score corrected based on the score difference to further promote the prediction error to approximate the true score.

$$\hat{s}_{\theta,t}(\hat{X}_t) = \hat{s}_{\theta,t}(\hat{X}_t)/\omega \quad (24)$$

We theoretically prove that the score difference helps to correct the score.

D DETAILS FOR EXPERIMENT

We provide detailed parameters for experiments related to §5, as shown in Tables 10 and 11. In particular, we differentiate the relevant parameters for sampling methods with correctors and those without correctors.

E SAMPLING ALGORITHM

In this subsection, we present the sampling algorithm procedure for S++, as shown in algorithm1. Additionally, our method can be naturally integrated into the reverse sampling of various diffusion models, greatly improving the generation quality of sampling methods without corrector. For methods with corrector, we can significantly reduce the correction time interval by introducing a truncation time step. Specifically, we only apply the corrector when the time step exceeds t_c further reducing computational cost.

Algorithm 1 The S++ sampling algorithm.

Input: pretrained real diffusion model $s_{\theta,t}(\cdot)$; Trained fake diffusion model $s_{\psi,t}(\cdot)$; Correction step size λ ; Cut-off time t_c
Initialize: $X_N \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
for $j = 1$ **to** M **do**
 $s_{\theta,N}(X_N) \leftarrow (s_{\theta,N}(X_N) + \lambda_1(s_{\theta,N}(X_N) - s_{\psi,N}(X_N))) / \omega_1$
 $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 $X_N \leftarrow X_N + \epsilon_N s_{\theta,N}(X_N) + \sqrt{2\epsilon_N} z$
end for
for $i = N - 1$ **to** 0 **do**
 $s_{\theta,i}(X_i) \leftarrow (s_{\theta,i}(X_i) + \lambda_2(s_{\theta,i}(X_i) - s_{\psi,i}(X_i))) / \omega_2$
 $X'_{i-1} \leftarrow (2 - \sqrt{1 - \beta_i})X_i + \beta_i s_{\theta,i}(X_i)$
 $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 $X_i \leftarrow X'_{i-1} + \sqrt{\beta_{i+1}} z$
if $i \leq t_c$ **then**
for $j = 1$ **to** M' **do**
 $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 $X_i \leftarrow X_i + \beta_i \epsilon_i s_{\theta,i+1}(X_{i+1}) + \sqrt{2\epsilon_i} z$
end for
end if
end for
return X_0

Model	Hyper.	Comm.	Enzymes	Grid	QM9	ZINC250k
GDSS-OC-S++	M	400	420	350	400	400
	λ_1	0.2	0.0008	0.06	1.19	2.5
	ω_1	0.998	1.0	1.0	1.09	1.0
	λ_2	0	0	0	0	0
	ω_2	1.0	1.0	1.0	1.0	1.0
HGDM-OC-S++	M	280	310	280	240	220
	λ_1	0.02	0.0	0.02	0.1	0.025
	ω_1	1.0	1.0	1.0	1.0	1.07
	λ_2	0.36	0.0	0.0	0.36	0.0
	ω_2	1.0	1.0	1.0	0.78	1.0
GSDM-OC-S++	M	200	400	400	-	-
	λ_1	0.0	0.0	0.0	-	-
	ω_1	1.0	1.0	1.0	-	-
	λ_2	0.0	0.0	0.0	-	-
	ω_2	1.0	1.0	1.0	-	-

Table 10: Experimental parameters for sampling methods with and without a corrector (OC).

F ADDITIONAL EXPERIMENTS

To demonstrate the superiority of S++, we selected generative models other than diffusion models as baseline models for comparison. GraphVAE (Simonovsky & Komodakis, 2018) is a graph generation model based on variational autoencoders; DeepGMG (Li et al., 2018) is a deep generative model that generates graphs in a sequential, pnode-by-node manner; GraphAF (Shi et al., 2020) is an autoregressive flow-based model. GraphRNN (You et al., 2018) is an autoregressive model using recurrent neural networks to generate graphs; EDP-GNN (Niu et al., 2020) is a score-based generative model using energy-based dynamics. GraphEBM (Liu et al., 2021) is an energy-based generative model that generates molecules by minimizing energy through Langevin dynamics, which is categorized as a one-shot generative method. We provide detailed comparative experiments in Tables 12 and 13, and the results show that our method significantly outperforms the baseline models and other generative models.

Model	Hyper.	Comm.	Enzymes	Grid	QM9	ZINC250k
GDSS-WC-S++	M	400	420	350	400	400
	λ_1	0.2	0.0008	0.06	1.19	2.5
	ω_1	0.998	1.0	1.0	1.09	1.0
	t_c	0.2	0.45	0.055	0.1	0.4
	λ_2	0	0	0	0	0
	ω_2	1.0	1.0	1.0	1.0	1.0
HGDM-WC-S++	M	280	200	360	240	220
	λ_1	0.02	0.0	0.18	0.1	0.25
	ω_1	1.0	1.0	1.0	1.0	1.07
	t_c	0.2	0.5	0.1	0.65	0.6
	λ_2	0.36	0.0	0.0	0.0	0.0
	ω_2	1.0	1.0	1.0	1.44	0.87
GSDM-WC-S++	M	200	400	400	-	-
	λ_1	0.0	0.0	0.0	-	-
	ω_1	1.0	1.0	1.0	-	-
	t_c	0.05	0.70	0.45	-	-
	λ_2	0.0	0.0	0.0	-	-
	ω_2	1.0	1.0	1.0	-	-

Table 11: Experimental parameters for sampling methods with and without a corrector (WC).

Dataset Info.	Community-small Synthetic, $12 \leq V \leq 20$				Enzymes Real, $10 \leq V \leq 125$				Grid Synthetic, $100 \leq V \leq 400$			
	Deg.	Clus.	Orbit	Avg.	Deg.	Clus.	Orbit	Avg.	Deg.	Clus.	Orbit	Avg.
DeepGMG	0.220	0.950	0.400	0.523	-	-	-	-	-	-	-	-
GraphRNN	0.080	0.120	0.040	0.080	0.017	0.062	0.046	0.042	0.064	0.043	0.021	0.043
GraphAF	0.18	0.20	0.02	0.133	1.669	1.283	0.266	1.073	-	-	-	-
GraphDF	0.06	0.12	0.03	0.070	1.503	1.061	0.202	0.922	-	-	-	-
GraphVAE	0.350	0.980	0.540	0.623	1.369	0.629	0.191	0.730	1.619	0.0	0.919	0.846
EDP-GNN	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-OC	0.050	0.132	0.011	0.064	0.052	0.627	0.249	0.309	0.270	0.009	0.034	0.070
GDSS-OC-S++	0.021	0.061	0.005	0.029	0.067	0.099	0.007	0.058	0.105	0.004	0.061	0.066
GDSS-WC	0.045	0.088	0.007	0.045	0.044	0.069	0.002	0.038	0.111	0.005	0.070	0.070
GDSS-WC-S++	0.019	0.062	0.004	0.028	0.031	0.050	0.003	0.028	0.105	0.004	0.061	0.057
HGDM-OC	0.065	0.119	0.024	0.069	0.125	0.625	0.371	0.374	0.181	0.019	0.112	0.104
HGDM-OC-S++	0.021	0.034	0.005	0.020	0.080	0.500	0.225	0.268	0.023	0.034	0.004	0.020
HGDM-WC	0.017	0.050	0.005	0.024	0.045	0.049	0.003	0.035	0.137	0.004	0.048	0.069
HGDM-WC-S++	0.021	0.024	0.004	0.016	0.040	0.041	0.005	0.029	0.123	0.003	0.047	0.058
GSDM-OC	0.142	0.230	0.043	0.138	0.930	0.867	0.168	0.655	1.996	0.0	1.013	1.003
GSDM-OC-S++	0.011	0.016	0.001	0.009	0.012	0.087	0.011	0.037	1.2e-4	0.0	1.2e-4	0.066
GSDM-WC	0.011	0.016	0.001	0.009	0.013	0.088	0.013	0.038	0.002	0.0	0	7.2e-5
GSDM-WC-S++	0.011	0.016	0.001	0.009	0.011	0.086	0.010	0.036	5.0e-5	0.0	1.1e-5	0.066

Table 12: Additional experiments on generic graph datasets.

Method	QM9			ZINC250k		
	Val. w/o corr. (%) \uparrow	NSPDK MMD \downarrow	FCD \downarrow	Val. w/o corr. (%) \uparrow	NSPDK MMD \downarrow	FCD \downarrow
GraphAF	67.00	0.020	5.268	68.00	0.044	16.289
GraphDF	82.67	0.063	10.816	89.03	0.176	34.202
MoFlow	91.36	0.017	4.467	63.11	0.046	20.931
EDP-GNN	47.52	0.005	2.680	82.97	0.049	16.737
GraphEBM	8.22	0.030	6.143	5.29	0.212	35.471
GDSS-OC	73.49	0.015	4.584	41.84	0.047	20.53
GDSS-OC-S++	93.74	0.001	1.661	59.50	0.050	16.79
GDSS-WC	94.91	0.004	2.550	95.83	0.019	14.66
GDSS-WC-S++	93.79	0.001	1.661	93.15	0.012	12.70
HGDM-OC	92.22	0.005	3.164	66.47	0.033	21.38
HGDM-OC-S++	94.95	0.003	2.512	67.12	0.034	20.79
HGDM-WC	98.02	0.002	2.147	93.26	0.016	17.69
HGDM-WC-S++	97.03	0.001	2.001	91.03	0.016	16.24

Table 13: Additional experiments on QM9 and ZINC250k datasets.

G INSENSITIVITY OF λ

We emphasize that S++ is insensitive to λ , since performance gain can always be achieved over a wide range of λ , as shown in Tables 14 and 15.

Table 14: FCD(\downarrow) on GDSS-OC baseline and QM9 under different parameters λ ($\lambda = 0$ represents the baseline)

λ	0	1.17	1.18	1.19	1.20	1.21
FCD	4.583	1.768	1.756	1.754	1.761	1.762

Table 15: Degree(\downarrow) on GDSS-OC baseline and Community-small under different parameters λ ($\lambda = 0$ represents the baseline)

λ	0	0.18	0.19	0.20	0.21	0.22
Degree	0.05	0.023	0.024	0.022	0.024	0.026

H TWO BIAS IN IMAGE AND GRAPH

Considering that there have been relevant studies on starting bias and exposure bias in the image field, this section will explain the differences between the two biases in the image and graph fields from three aspects: data scale, data structure, and network performance.

1-a Starting Bias in Images. DPM-Fixes (Lin et al., 2024) first discovered that conventional noise scheduling strategies cannot guarantee that the maximum forward perturbation distribution follows a standard Gaussian distribution:

$$\mathbf{x}_T = 0.068265\mathbf{x}_0 + 0.997667\epsilon_T \quad (25)$$

This shows a slight deviation from the reverse sampling starting point of standard Gaussian. DPM-Fixes forces forward \mathbf{x}_T to follow standard Gaussian by adjusting noise scheduling scale, benefiting from sufficient image data scale and strong network performance, enabling accurate noise (or score) prediction in high-noise states. DPM-Leak (Everaert et al., 2024) estimates the maximum forward perturbation distribution during training based on pixel modeling and uses it as new sampling starting points, benefiting from image data structure and scale that allows pixels to be independent and follow Gaussian distribution.

1-b Starting Bias in Graphs. Unlike image data, limited by data scale and network performance, models struggle to accurately predict noise (or score) from high-noise states. Therefore, the baseline model adopts a conservative strategy during training, resulting in the maximum forward perturbation distribution falling far short of standard Gaussian. While this avoids high-noise states, it introduces significant starting bias. Consequently, strategies like DPM-Fixes that force training \mathbf{x}_T to follow standard Gaussian are unsuitable; due to node and edge interdependencies in graph data and high sparsity characteristics, we cannot simply assume nodes or edges follow Gaussian distribution to estimate maximum forward perturbation distribution, making DPM-Leak unsuitable as well.

2-a Exposure Bias in Images. In images, exposure bias refers to the mismatch between forward process \mathbf{x}_t and reverse process $\hat{\mathbf{x}}_t$, with differences accumulating throughout sampling, ultimately affecting generation quality. Many current image exposure bias works assume no starting bias exists, focusing on sampling process bias, as starting bias in images is indeed quite minimal.

2-b Exposure Bias in Graphs. Unlike the minor signal leakage in images, graph diffusion models have significant starting bias, meaning severe exposure bias exists after the first sampling step. In other words, graph exposure bias isn’t solely caused by network prediction errors and sampling iteration accumulation but is severely impacted by starting bias. Therefore, addressing graph exposure bias requires first resolving starting bias.

In conclusion, we emphasize that starting bias in graph diffusion models is a more acute and unique problem. Although exposure bias exists in graphs, image-based solutions cannot be simply imitated. We are the first work focusing on starting bias in graph diffusion models, proposing a simple yet effective solution. We aim to draw attention from relevant researchers, hoping they consider bias analysis while developing graph diffusion models.

I S++ AND EXISTING SOLUTIONS IN IMAGES

Although there are many solutions to solve the exposure bias in current images, these solutions cannot replace S++. In this section, we choose to compare with similar works (Li et al., 2024; Ning et al., 2024) to S++ in detail because they are plug-and-play solutions that do not introduce new components.

TS-DPM(Li et al., 2024) proposes searching an optimal timestep s during sampling. TS-DPM relies on two fundamental assumptions: (a) image pixels are independent and follow Gaussian distribution - Eq.13 in Appendix J of [1]; (b) sample pixel variance approximates population variance - Eq.20 in Appendix J of TS-DPM. These assumptions are based on large image datasets and large number of pixels. However, nodes and edges in graphs are highly sparse. Specifically, many graph datasets are small (Community-small, Enzymes, and Grid have fewer than 1000 samples). Both assumptions do not hold for graph data, making TS-DPM inapplicable to graph diffusion models.

ADM-ES(Ning et al., 2024) proposes reducing $s_{\theta,t}(\cdot)$ (originally $\epsilon_{\theta,t}(\cdot)$) during sampling to mitigate exposure bias. However, this approach does not involve the angle of the vector $s_{\theta,t}(\cdot)$. Our approach addresses this limitation:

$$s_{\theta,t}(X_t) = (s_{\theta,t}(X_t) + \lambda(s_{\theta,t}(X_t) - s_{\psi,t}(X_t)))/\omega \quad (26)$$

when $\lambda = 0$, $s_{\theta,t}(X_t) \leftarrow s_{\theta,t}(X_t)/\omega$, which is equivalent to work ADM-ES. In other words, ADM-ES is a special case of S++, while $(s_{\theta,t}(X_t) - s_{\psi,t}(X_t))$ provides angle information.

For fair comparison with ADM-ES, we introduce $\lambda(s_{\theta,t}(X_t) - s_{\psi,t}(X_t))$ for each magnitude factor ω , with λ uniformly set to 0.5, to examine whether $(s_{\theta,t}(X_t) - s_{\psi,t}(X_t))$ brings improvements over ADM-ES. Tables 16 and 17 demonstrate that the angle information in S++ leads to significant gains. We conducted experiments between S++ and ADM-ES as following:

Table 16: FCD(\downarrow) on QM9 without corrector.

ω	0.7	0.8	0.9	1.0	1.1
GDSS-OC-ES	3.57	3.417	3.768	4.584	5.517
GDSS-WC-S++	2.94	2.814	3.198	4.187	5.201

Table 17: FCD(\downarrow) on QM9 with corrector.

ω	0.9	1.0	1.1	1.2	1.3
GDSS-WC-ES	2.809	2.552	2.319	2.301	2.542
GDSS-WC-S++	2.321	2.034	1.858	1.852	2.152