# NVDiff: Graph Generation through the Diffusion of Node Vectors

Xiaohui Chen[1], Yukun Li[1], Aonan Zhang[2], Li-Ping Liu[1*]

[1]Department of Computer Science, Tufts University
[2]ByteDance Inc.

## Abstract

Learning to generate graphs is challenging as a graph is a set of pairwise connected, unordered nodes encoding complex combinatorial structures. Recently, several works have proposed graph generative models based on normalizing flows or score-based diffusion models. However, these models need to generate nodes and edges in parallel from the same process, whose dimensionality is unnecessarily high. We propose NVDiff, which takes the VGAE structure and uses a score-based generative model (SGM) as a flexible prior to sample node vectors. By modeling only node vectors in the latent space, NVDiff significantly reduces the dimension of the diffusion process and thus improves sampling speed. Built on the NVDiff framework, we introduce an attention-based score network capable of capturing both local and global contexts of graphs. Experiments indicate that NVDiff significantly reduces computations and can model much larger graphs than competing methods. At the same time, it achieves superior or competitive performances over various datasets compared to previous methods.

## 1 Introduction

Graph generation learns a distribution over graphs that helps summarize, predict, and explore various domains such as molecules and social networks (Jin, Barzilay, and Jaakkola 2018; Watts and Strogatz 1998). This task is challenging because it requires a generative model to learn complex combinatorial patterns hidden in the graph data. One category of neural generative models sequentially generates nodes and edges of a graph (You et al. 2018; Liao et al. 2019). However, the generation probability of such a model varies when node ordering changes. Fitting such a model often requires node orderings, which often do not exist in training data. Furthermore, it is hard to train such models by maximizing the data likelihood (Chen et al. 2021).

This paper focuses on another category of graph generative models, which assumes nodes are exchangeable (Aldous 1981; Hoover 1979). Those models generate graphs as adjacency matrices invariant in probability when jointly shuffling rows and columns with a permutation $\pi$. There are two benefits to using exchangeable graphs. First, since node ordering does not affect the generation probability, calculating the

likelihood of the graph only requires one pass with an arbitrary node ordering. Second, graph generation can be done in one-shot, i.e., parallel rather than sequential.

One early example of deep exchangeable graph generative models is VGAE (Kipf and Welling 2016). Based on the variational auto-encoder (VAE) framework (Kingma and Welling 2013), VGAE represents the graph by latent node vectors. The priors of those latent vectors follow a multivariate Gaussian distribution, and they use an encoder to infer the variational posterior. VGAE uses an over-simplified prior and decoder. Its node vectors are independent, thus they cannot "coordinate" to play different roles in a meaningful graph.

Recently, normalizing flows (NF) (Rezende and Mohamed 2015; Liu et al. 2019; Lippe and Gavves 2020; Xu et al. 2021) are introduced for graph generation by learning a reversible mapping from a Gaussian distribution to a continuous latent graph representation. However, the reversibility assumption restricts model design for NF models. For instance, graph normalizing flows (GNF) (Liu et al. 2019) detaches NF with Graph Auto-Encoders, thus not guaranteed to maximize the overall likelihood. GraphCNF (Lippe and Gavves 2020) incorporates both node and edge latent vectors in NF and thus has a generative space with large dimensionality.

Score-based generative models (SGMs) introduce a diffusion process that maps observations into a noise distribution—reversing the diffusion process results in a generative model, which can be learned by score matching using, e.g., noise conditional score networks (NCSN) (Song and Ermon 2019). SGMs demonstrate superior performance than NF in generating realistic images and are recently adapted for one-shot graph generation. Niu et al. (2020) proposed EDP-GNN to model graphs using discrete-time SGMs. Jo, Lee, and Hwang (2022) introduced GDSS, a continuous-time graph diffusion process for nodes and edges. Learning the graph generative process requires solving a system of reverse-time SDEs (Anderson 1982; Song et al. 2020). Both EDP-GNN and GDSS directly apply diffusion processes on edges. These methods are not scalable since their inference requires reversing a $\mathcal{O}(N^2)$-dimensional diffusion process. GDSS ignores the hierarchical interpretation between nodes and edges by modeling node and edge representations in parallel.

In observation of those limitations, we propose NVDiff, a novel graph generative model built on latent score-based generative models (LSGM) (Vahdat, Kreis, and Kautz 2021).

---
[*]Corresponding email: liping.liu@tufts.edu

The generative process first samples a latent representation of the graph and then decodes nodes and edges conditioning on the latent representation. In principle, the dimension of the latent space is arbitrary. To guarantee node exchangeability, we propose a diffusion process on latent node vectors, thus significantly reducing the dimension of the diffusion process from $\mathcal{O}(N^2)$ to $\mathcal{O}(N)$ and boosting sampling speed. Built on the NVDiff framework, we introduce an attention-based score network capturing local and global contexts within the graph by leveraging the self-attention mechanism with a contextual vector. NVDiff optimizes the variational lower bound of log-likelihood, which can be optimized end-to-end. Our empirical study demonstrates that NVDiff outperforms previous methods over various graph generation tasks. By diffusing only the node vectors, NVDiff can learn to generate larger graphs with high quality in one-shot.

## 2 Background: Score-based Generative Models

Denote a continuous time diffusion process as $(\mathbf{Z}^t)_{t \in [0,1]}$ with base distribution $q_0(\mathbf{Z}^0)$. In practice, $q_0$ can either be the data distribution or latent variable distribution. Assuming the diffusion process follows an Itô SDE:

$$d\mathbf{Z}^t = f(t)\mathbf{Z}^t dt + g(t)d\mathbf{W}^t, \tag{1}$$

where $\mathbf{W}^t$ is the standard Wiener process. $f(t)$ and $g(t)$ are scalar functions denoting drift and diffusion coefficients, which are chosen such that $q_1(\mathbf{Z}^1) = \mathcal{N}(\mathbf{Z}^1; \mathbf{0}, \mathbf{I})$. The forward SDE gradually corrupt the base distribution to a random noise distribution that is easy to sample from, e.g., Gaussian distribution. To sample from the base distribution, we first draw $\mathbf{Z}^1$ from $\mathcal{N}(\mathbf{Z}^1; \mathbf{0}, \mathbf{I})$ and run a reverse-time SDE:

$$d\mathbf{Z}^t = [f(t)\mathbf{Z}^t - g(t)^2 \nabla_{\mathbf{Z}^t} \log q_t(\mathbf{Z}^t)]dt + g(t)d\bar{\mathbf{W}}^t, \tag{2}$$

where $\bar{\mathbf{W}}^t$ is the reverse-time standard Wiener process, and $q_t(\mathbf{Z}^t)$ is the marginal distribution of the forward SDE at time $t$. Since $q_t(\mathbf{Z}^t)$ is analytically intractable, we use a noise conditional score network $\epsilon_\theta(\mathbf{Z}^t, t)$ to fit the score function as:

$$\min_\theta \ \mathbb{E}_{t \sim U[0,1]} \Big[ \frac{g(t)^2}{2} \mathbb{E}_{q_0(\mathbf{Z}^0)} \mathbb{E}_{q(\mathbf{Z}^t|\mathbf{Z}^0)} [\| \nabla_{\mathbf{Z}^t} \log q_t(\mathbf{Z}^t)$$
$$- \epsilon_\theta(\mathbf{Z}^t, t)\|_2^2] \Big], \tag{3}$$

where $q(\mathbf{Z}^t|\mathbf{Z}^0)$ denotes the transition kernel of forward SDE. Replacing $\nabla_{\mathbf{Z}^t} \log q_t(\mathbf{Z}^t)$ with the score network $\epsilon_\theta(\mathbf{Z}^t, t)$ in Eqn. (2) derives a score-based generative model (SGM). We denote the marginal distribution of SGM at $t = 0$ as $p_\theta(\mathbf{Z})$.

## 3 NVDiff: Graph Generation through Node Variable Diffusion

### 3.1 Graphs as Exchangeable Adjacency Matrices

We consider an attributed graph $G$ with $N$ nodes and denote it by $(\mathbf{A}, \mathbf{X})$. Here $\mathbf{A} \in \mathbb{R}^{N \times N \times K^e}$ denotes edge features, $\mathbf{X} \in \mathbb{R}^{N \times K^v}$ denotes node features, and feature dimensions are $K^e$ and $K^v$. We mainly consider the case where nodes

and edges are categorical data with one-hot vector representation.

An exchangeable generative model $p$ requires $p(\mathbf{A}, \mathbf{X}) \stackrel{d}{=} p(\mathbf{A}[\pi, \pi, :], \mathbf{X}[\pi, :])$ for any permutation $\pi : [N] \to [N]$. $\mathbf{A}[\pi, \pi, :]$ represents the resulting matrix by jointly permuting the first and the second coordinates of $\mathbf{A}$, and $\mathbf{X}[\pi, :]$ represents permuting rows of $\mathbf{X}$ with $\pi$. With proper model design in the variational auto-encoder framework (VAE) by introducing intermediate latent variables, we can build flexible graph generative models without sacrificing exchangeability.

### 3.2 The VAE Framework

Variational Graph Auto-Encoder (VGAE) (Kipf and Welling 2016) is a variational auto-encoder that models the distribution of adjacency matrices. It imposes a Gaussian prior $p(\mathbf{Z})$ over node vectors. During training, it maps $(\mathbf{A}, \mathbf{X})$ to $\mathbf{Z}$ using an encoder $q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X})$ and then recover $\mathbf{A}$ using a decoder $p_\psi(\mathbf{A}|\mathbf{Z})$.

Our work implements the VGAE framework and makes several important differences. (1) We use SGM as our prior $p_\theta(\mathbf{Z})$ to increase expressiveness. (2) We use decoder $p_\psi(\mathbf{A}, \mathbf{X}|\mathbf{Z})$ to model both nodes and edges. (3) We make adjustment to the encoder $q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X})$ to facilitate model training. We highlight our framework in Figure 1.

Learning VAE requires maximizing the variational lower bound $\mathcal{L}(\phi, \psi, \theta)$ of log-likelihood $\log p(\mathbf{A}, \mathbf{X})$ as

$$\max_{\phi, \psi, \theta} \ \mathcal{L}(\phi, \psi, \theta) = \mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{A},\mathbf{X})} \big[ \log p_\psi(\mathbf{A}, \mathbf{X}|\mathbf{Z}) \big]$$
$$+ \mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{A},\mathbf{X})} \big[ \log p_\theta(\mathbf{Z}) \big] + \mathbb{H} \big[ q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X}) \big]. \tag{4}$$

In principle, $\mathbf{Z}$ can be a latent variable with any dimension. Moderately compressed $\mathbf{Z}$ enjoys compact representation and improves sampling speed of SGM. In NVDiff, we choose $\mathbf{Z} \in \mathbb{R}^{N \times d}$, where each row corresponds to a latent node vector. To generate graphs, we first draw node vectors $\mathbf{Z} \sim p_\theta(\mathbf{Z})$ and draw $\mathbf{A}, \mathbf{X} \sim p_\psi(\mathbf{A}, \mathbf{X}|\mathbf{Z})$. Edge information is encoded in mutually-dependent node vectors through a learnable SGM prior $p_\theta(\mathbf{Z})$.

### 3.3 Attention-based Score Network with Contextual Embedding

In NVDiff, learning $p_\theta(\mathbf{Z})$ is equivalent to learning the score network $\epsilon_\theta(\mathbf{Z}^t, t) : \mathbb{R}^{N \times d} \times \mathbb{R} \to \mathbb{R}^{N \times d}$. To guarantee exchangeability of the SGM prior, $\epsilon_\theta(\mathbf{Z}^t, t)$ needs to be permutation-equivariant in rows, i.e., $\epsilon_\theta(\mathbf{Z}^t, t)[\pi, :] = \epsilon_\theta(\mathbf{Z}^t[\pi, :], t)$ for any permutation $\pi$ (Niu et al. 2020).

To guarantee permutation-equivariance without losing expressiveness, we use multi-head self-attention layers (MHA) to construct $\epsilon_\theta(\mathbf{Z}^t, t)$. Following Song et al. (2020), we transform $t$ to its positional encoding vector $\mathbf{t} = \mathrm{pe}(t)$ as an input to the initial mutli-layer perceptron (MLP) layer:

$$\mathbf{H}_0 = \mathrm{mlp}^{\mathrm{in}}([\mathbf{Z}^t; \mathbf{1}\mathbf{t}^\top]), \tag{5}$$

where $[\cdot; \cdot]$ represents stacking two matrices by aligning their rows, and the $\mathrm{mlp}^{\mathrm{in}}(\cdot)$ applies to rows of its argument.

Then we stack $L$ MHA layers to process the input $\mathbf{H}_0$ and incorporate a learnable global vector $\mathbf{g}_0$ to carry graph
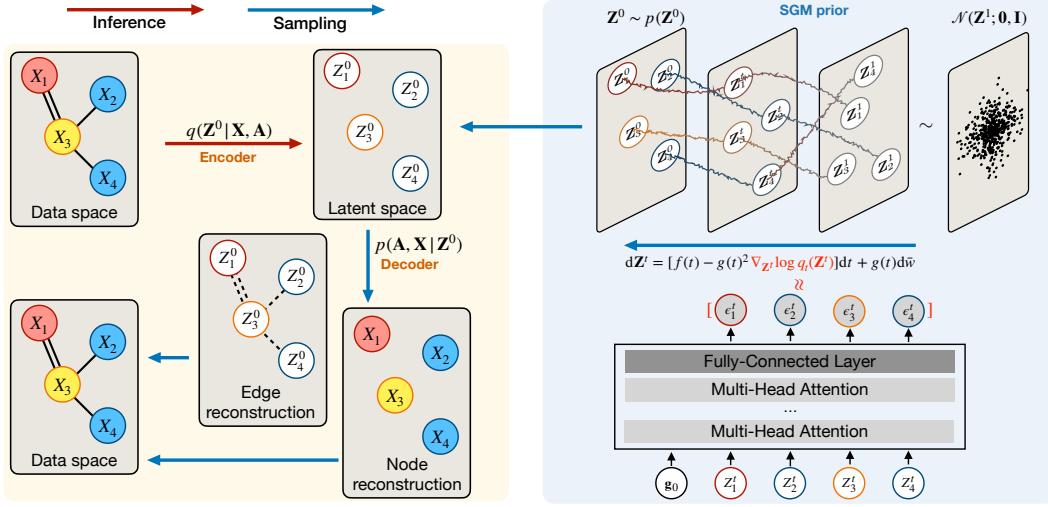
Figure 1: In NVDiff, a graph $(\mathbf{A}, \mathbf{X})$ is mapped to a set of latent node vectors $\mathbf{Z}^0$ through an encoder, the joint distribution of the latent node vectors is model by a diffusion prior. To generate a graph, NVDiff first samples $\mathbf{Z}^0$ by reversing the diffusion process, then generates nodes and edges by decoding the latent representation.

information through MHA layers:

$$(\mathbf{g}'_\ell, \mathbf{H}_\ell) = \mathrm{mha}([\mathbf{g}_{\ell-1}; \mathbf{H}_{\ell-1}]),$$
$$\mathbf{g}_\ell = \mathbf{g}_{\ell-1} + \mathbf{g}'_\ell, \quad \ell = 1, \dots, L. \tag{6}$$

Finally, we predict the score from the representation $\mathbf{H}_L$:

$$\hat{\epsilon} = \mathrm{mlp}^{\mathrm{out}}(\mathbf{H}_L). \tag{7}$$

The input of the score network $\mathbf{Z}^t$ is sampled from a transition distribution $q(\mathbf{Z}^t|\mathbf{Z}^0) = \mathcal{N}(\mathbf{Z}^t; \boldsymbol{\mu}_t(\mathbf{Z}^0), \sigma_t \mathbf{I})$, which has a closed form with $\boldsymbol{\mu}_t(\cdot)$ and $\sigma_t$ determined by $t$. Since $\mathbf{Z}^t$ naturally carries blurred structural information of each node, we use a contextual vector $\mathbf{g}_0$ to carry the global information by attending all the node vectors at the first MHA layer. The following MHA layers gradually mix information between the contextual vector $\mathbf{g}$ and the node vectors $\mathbf{H}$ through self-attention until getting the output score.

To draw a sample $\mathbf{Z}^0$ from $p_\theta(\mathbf{Z})$, we first draw $\mathbf{Z}^1$ from $\mathcal{N}(\mathbf{Z}^1; \mathbf{0}, \mathbf{I})$ and run the reverse-time SDE (Eqn. (2)) with true score function $\nabla_{\mathbf{Z}^t} \log q_t(\mathbf{Z}^t)$ replaced by the score network $\epsilon_\theta(\mathbf{Z}^t, t)$ (Song et al. 2020).

### 3.4 Graph Decoding

The decoder predicts node attributes, edges, and edge attributes from node representations $\mathbf{Z}$. We first transform $\mathbf{Z}$ to get the node representation $\mathbf{S} \in \mathbb{R}^{N \times d'}$ and the edge representation $\mathbf{R} \in \mathbb{R}^{N \times N \times d''}$ from the penultimate layer of the decoder.

$$(\mathbf{S}, \mathbf{R}) = \mathrm{gnn}_\psi(\mathbf{Z}, \mathbf{E}), \ \mathbf{E} = [(\mathbf{z}_i - \mathbf{z}_j)^2 : i, j = 1, \dots, N]. \tag{8}$$

Here $\mathbf{E} \in \mathbb{R}^{N \times N \times d}$ is a tensor, whose $(i, j)$-th vector slice denotes the element-wise squared difference between two node vectors $\mathbf{z}_i$ and $\mathbf{z}_j$. $\mathrm{gnn}_\psi(\cdot, \cdot)$ is a $M$-layer graph neural network that computes on a fully-connected graph with $\mathbf{Z}, \mathbf{E}$

being the node and edge inputs. The $\ell$-th layer of the GNN updates nodes and edges from $(\mathbf{S}^{\ell-1}, \mathbf{R}^{\ell-1})$ to $(\mathbf{S}^\ell, \mathbf{R}^\ell)$ as follow:

$$\mathbf{r}^\ell_{i,j} = \mathrm{mlp}^e([\mathbf{W}^e \mathbf{s}^{\ell-1}_i + \mathbf{W}^e \mathbf{s}^{\ell-1}_j; \mathbf{r}^{\ell-1}_{i,j}]),$$
$$\mathbf{s}^\ell_i = \mathrm{mlp}^v\Big(\Big[\sum_{j:j\neq i}^N \mathbf{W}^v \mathbf{r}^\ell_{i,j}; \mathbf{s}^{\ell-1}_i\Big]\Big). \tag{9}$$

Here $\mathbf{s}^\ell_i$ is the $i$-th row in $\mathbf{S}^\ell$, $\mathbf{r}^\ell_{i,j}$ is the $(i, j)$-th vector slice in $\mathbf{R}^\ell$, we denote $(\mathbf{S}^0, \mathbf{R}^0) := (\mathbf{Z}, \mathbf{E})$ as inputs to the GNN, and $(\mathbf{S}, \mathbf{R}) := (\mathbf{S}^M, \mathbf{R}^M)$. We recover node features and edge features respectively:

$$p(\mathbf{A}, \mathbf{X}|\mathbf{Z}) = p(\mathbf{A}|\mathbf{R})p(\mathbf{X}|\mathbf{S}) = \prod_{i=1}^N p(\mathbf{x}_i|\mathbf{s}_i) \prod_{j<i} p(\mathbf{a}_{i,j}|\mathbf{r}_{i,j}), \tag{10}$$

where $\mathbf{x}_i$ denotes the $i$-th row in $\mathbf{X}$ and $\mathbf{a}_{i,j}$ denotes the $(i, j)$-th vector slice in $\mathbf{A}$. Since the adjacency matrix is symmetric, we consider the lower triangle only.

Now we are ready to recover individual features for edges and nodes. When we consider an edge feature, we treat a non-edge ($\mathbf{a}_{i,j} = \mathbf{0}$) separately from an actual edge label. We specify $p(\mathbf{a}_{i,j}|\mathbf{r}_{i,j})$ with two set of probabilities:

$$p(\mathbf{a}_{i,j} = \mathbf{0}|\mathbf{r}_{i,j}) = \mathrm{mlp}^b_\psi(\mathbf{r}_{i,j}), \tag{11}$$
$$p(\mathbf{a}_{i,j}|\mathbf{a}_{i,j} \neq \mathbf{0}, \mathbf{r}_{i,j}) = \mathrm{Categorical}(\mathbf{a}_{i,j}; \mathrm{mlp}^e_\psi(\mathbf{r}_{i,j})).$$

Here $\mathrm{mlp}^b_\psi(\cdot)$ calculates the Bernoulli probability of the non-existence of edge $(i, j)$, and $\mathrm{mlp}^e_\psi(\cdot)$ calculates the categorical probabilities for the label of edge $(i, j)$.

Then we consider node features and define:

$$p(\mathbf{x}_i|\mathbf{s}_i) = \mathrm{Categorical}(\mathbf{x}_i; \mathrm{mlp}^n_\psi(\mathbf{s}_i)). \tag{12}$$

Here $\mathrm{mlp}_\psi^n(\cdot)$ is applied to get the probabilities of the categorical distribution. We use sigmoid activation function in the last layer of $\mathrm{mlp}_\psi^b(\cdot)$ and softmax activation function in the last layer of $\mathrm{mlp}_\psi^e(\cdot)$ and $\mathrm{mlp}_\psi^n(\cdot)$.

## 3.5 Graph Encoding

We use a mean-field distribution as the encoder and parameterize it with a GNN:

$$q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X}) = \mathcal{N}(\mathbf{Z}; \mathbf{M}, \sigma^2\mathbf{I}), \ \mathbf{M} = \mathrm{gnn}_\phi(\mathbf{A}, \mathbf{X}). \quad (13)$$

Here the encoder is a Gaussian distribution, whose mean is computed from a GNN. The architecture details are described in Appendix. To simplify the optimization objective, we fix the variance of the encoder and only optimize its mean parameters. By doing so, we constrain flexibility of the encoder and encourage the SGM prior to match the aggregated posterior. When we apply a GNN to $(\mathbf{A}, \mathbf{X})$, we treat non-edges as an extra edge-type and also pass messages through them. Following Satorras, Hoogeboom, and Welling (2021), we also inject random noise to each node to break the symmetry of the graph such that nodes are separable. We find these measures improve the model performance.

## 3.6 Training

In this section, we present the training procedure that iteratively optimizes the encoder, the decoder, and the diffusion prior. The training procedure is shown in Alg. 1. In the optimization objective $\mathcal{L}(\phi, \psi, \theta)$ in Eqn. (4), the first term $\mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X})}\big[\log p_\psi(\mathbf{A}, \mathbf{X}|\mathbf{Z})\big]$ is the reconstruction term. We apply the reparameterization technique (Kingma and Welling 2013) and use Monte Carlo samples to estimate the expectation. Vahdat, Kreis, and Kautz (2021) shows that the cross-entropy term can be calculated through score matching:

$$\mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X})}\big[\log p_\theta(\mathbf{Z})\big]$$
$$= \mathbb{E}_{t, \epsilon, q_\phi(\mathbf{Z}^0, \mathbf{Z}^t|\mathbf{A}, \mathbf{X})}\left[\frac{g(t)^2}{2}||\epsilon - \epsilon_\theta(\mathbf{Z}^t, t)||_2^2\right] \quad (14)$$

where $t \sim U[0, 1]$, $\epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$, and $q_\phi(\mathbf{Z}^0, \mathbf{Z}^t|\mathbf{A}, \mathbf{X}) = q(\mathbf{Z}^t|\mathbf{Z}^0)q_\phi(\mathbf{Z}^0|\mathbf{A}, \mathbf{X})$. Here $q(\mathbf{Z}^t|\mathbf{Z}^0)$ is a Gaussian transition kernel. The cross-entropy term is also estimated through Monte Carlo method. We describe the details of computing the cross-entropy term in Appendix. The entropy term $\mathbb{H}[q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X})]$ in Eqn. (4) is a constant because we fix the variance of the posterior distribution.

Recall that $\mathcal{L}(\phi, \psi, \theta) = \mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X})}\big[\log p_\psi(\mathbf{A}, \mathbf{X}|\mathbf{Z})\big] + \mathrm{KL}(q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X})||p_\theta(\mathbf{Z}))$. We use KL annealing to dynamically balance optimization between the KL term and the rest objective functions during training. KL annealing effectively prevents the variational posterior from converging to a premature diffusion prior.

Sampling novel and unique graphs requires the decoder to be robust to noise in node vectors, since the diffusion prior may generate node vectors unseen during training. We introduce a consistency regularization term to encourage the decoder to generalize better, by adapting to a perturbed latent

---

**Algorithm 1: Training Procedure for NVDiff**

**Input:** Graph dataset $\mathcal{G} = \{(\mathbf{A}_1, \mathbf{X}_1), \ldots\}$, encoder $q_\phi$, encoder $p_\psi$, diffusion prior $p_\theta$, transition kernel $q$, diffusion coefficient $g(\cdot)$
**for** $(\mathbf{A}, \mathbf{X}) \in \mathcal{G}$ **do**
  **[Optimize the encoder and the decoder]**
    Draw $\mathbf{Z}^0 \sim q_\phi(\mathbf{Z}|\mathbf{X}, \mathbf{A}), t \sim U[0, 1], \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$, $\mathbf{Z}^t \sim q(\mathbf{Z}^t|\mathbf{Z}^0)$.
    Optimize $\mathcal{L}_{\mathrm{VAE}} = -\log p_\psi(\mathbf{A}, \mathbf{X}|\mathbf{Z}^0) + \frac{g(t)^2}{2}||\epsilon - \epsilon_\theta(\mathbf{Z}^t, t)||_2^2$ over $\phi$ and $\psi$ with one step.
  **[Optimize the SGM prior]**
    Optimize $\mathcal{L}_{\mathrm{SGM}} = \frac{g(t)}{2}||\epsilon - \epsilon_\theta(\mathbf{Z}^t, t)||_2^2$ over $\theta$ with one step.
**end for**

---

variable $\tilde{\mathbf{Z}}$. We define:

$$\mathcal{L}_{\mathrm{Reg}} = \mathbb{E}_{\tilde{\epsilon}, q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X})}\big[-\log p_\psi(\mathbf{A}, \mathbf{X}|\tilde{\mathbf{Z}})\big],$$
$$\text{with } \tilde{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_1\mathbf{I}), \ \mathbf{Z} \sim q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X}), \ \tilde{\mathbf{Z}} = \mathbf{Z} + \tilde{\epsilon}. \quad (15)$$

Here $\sigma_1$ is a hyperparameter. We replace the reconstruction term with $\mathcal{L}_{\mathrm{Reg}}$ and finetune the decoder.

## 4 Related Works

**Score-based generative models.** Score-based generative models (Song and Ermon 2019; Ho, Jain, and Abbeel 2020; Song et al. 2020) have shown effective at the generation of data in various domains (Ho, Jain, and Abbeel 2020; Chen et al. 2020; Niu et al. 2020). In SGMs, data are first diffused with gradually increasing noise, and the model learns to estimate the score at each noise scale. During generation, an SGM first randomly draws some noise from a simple distribution and then decreases the noise scales sequentially to obtain a sample. Instead of generating data directly, Mittal et al. (2021); Sinha et al. (2021); Wehenkel and Louppe (2021); Vahdat, Kreis, and Kautz (2021) utilize SGM in latent space and turn it into a flexible prior by combining the VAE framework (Kingma and Welling 2013). Unlike approaches that use more complex prior (Dai and Wipf 2019; Klushyn et al. 2019), SGM prior addresses the posterior mismatch problem by construction (Sinha et al. 2021). Combining SGM prior into VAE framework also enables faster sampling as well as hierarchically generating data.

**Graph generation models.** Current graph generative models sample graphs either autoregressively or in one-shot. Autoregressive graph generative models generate graphs by adding nodes and edges sequentially (You et al. 2018; Li, Zhang, and Liu 2018; Li et al. 2018; Jin, Barzilay, and Jaakkola 2018; Liao et al. 2019; Shi et al. 2020; Zang and Wang 2020; Bacciu, Micheli, and Podda 2020; Dai et al. 2020; Goyal, Jain, and Ranu 2020; Tran et al. 2020; Luo, Yan, and Ji 2021; Bengio et al. 2021). Training such models requires a particular node ordering from the graph to pin down the generation sequence, thus unable to preserve the model exchangeability (Chen et al. 2021). Besides, the autoregressive models usually fail to learn the global structure of a graph as they can

Table 1: Performance on molecule generation. Numbers in bold indicate that the methods are the best, and numbers underlined are the second best. Extended NSPDK, Uniqueness metric as well as standard deviations are in Appendix.

| | QM9 | | | ZINC250K | | |
|---|---|---|---|---|---|---|
| | Validity↑ | FCD↓ | Time (s)↓ | Validity↑ | FCD↓ | Time (s)↓ |
| GraphDF | 82.67 | 10.816 | 5.35e4 | 89.03 | 34.202 | 6.03e4 |
| MoFlow | 91.36 | 4.467 | **4.60** | 63.11 | 20.931 | **24.5** |
| GraphCNF | 95.00 | <u>1.629</u> | 59.48 | <u>95.26</u> | <u>14.786</u> | 1.04e2 |
| EDP-GNN | 47.52 | 2.680 | 4.40e3 | 82.97 | 16.737 | 9.09e3 |
| GDSS | <u>95.72</u> | 2.900 | 1.14e2 | **97.01** | 14.656 | 2.02e3 |
| NVDiff | **95.79** | **1.131** | <u>10.71</u> | 85.63 | **4.019** | <u>90.2</u> |

not capture the long-term dependency. The other type of models generate graph in one-shot (Simonovsky and Komodakis 2018; Liu et al. 2018; Madhawa et al. 2019; Kajino 2019; Bradshaw et al. 2019; Grover, Zweig, and Ermon 2019; Lim et al. 2020; Jin, Barzilay, and Jaakkola 2020; Guo et al. 2020; Lippe and Gavves 2020; Samanta et al. 2020; Niu et al. 2020; Garcia Satorras et al. 2021; Jo, Lee, and Hwang 2022). Unlike autoregressive models, exchangeability can be preserved through carefully designing the model architecture. However, those models need to generate nodes and edges in parallel and can not scale up to large graphs. On par with the one-shot graph generative models, VGAEs (Kipf and Welling 2016; Mehta, Duke, and Rai 2019; Liu et al. 2019; Li et al. 2020) adopt the variational inference framework to learning node representations that are useful for downstream tasks, such as unsupervised learning or node clustering. However, The merits of using such approaches in graph generation, such as the exchangeable property and the parallel sampling scheme, have not been fully exploited. Our work builds on the VGAE framework and utilizes a learnable prior based on SGM.

## 5 Experiments

We design a set of experiments to demonstrate that NVDiff can generate high-quality graphs with latent node vectors. Specifically, in Section 5.2 we quantitatively evaluate the sample quality and sampling speed for selected baseline models and NVDiff over various metrics. In Section 5.3 we investigate NVDiff from the perspective of the generation dynamics, effectiveness of the contextual embedding, and graph visualization. More experiment results can be found in the Appendix.

### 5.1 Experiment Setup

Here we briefly discuss the datasets, baselines, and evaluation metrics we use in the experiments. Further details can be found in Appendix.
**Datasets.** We use six datasets, including two molecular datasets and four generic graph datasets. For molecular datasets, we consider QM9 (Ramakrishnan et al. 2014) and ZINC250K (Irwin et al. 2012). For generic graph datasets, we consider Community and Ego and the small versions of them.
**Baselines.** We consider the following baselines: GraphRNN and GRAN (You et al. 2018; Liao et al. 2019) are autoregressive models. GraphDF (Luo, Yan, and Ji 2021) is an autoregressive flow-based model. Moflow, GraphCNF (Zang and

Wang 2020; Lippe and Gavves 2020) are flow-based models that generate graphs in one-shot. EDP-GNN and GDSS (Niu et al. 2020; Jo, Lee, and Hwang 2022) are one-shot score-based graph generative models. For molecule generation, we compare NVDiff with GraphDF, Moflow, GraphCNF, EDP-GNN, and GDSS. We compare generic graph generation with GraphRNN, GRAN, GraphCNF, EDP-GNN, and GDSS.
**Evaluation metrics.** We evaluate the generated molecules with the following metrics: Validity is the fraction of generated molecules with valid chemical valency. And Uniqueness is the fraction of unique molecules. We evaluate the Maximum Mean Discrepancy (MMD) (Gretton et al. 2012) between the test molecules and the generated molecules in terms of the Neighborhood subgraph pairwise distance kernel (NSPDK) (Costa and De Grave 2010), which considers both structure and labels of the graph. Fréchet ChemNet Distance (FCD) (Preuer et al. 2018) evaluates the diversity and the quality of the validly generated molecules by measuring the distance between feature vectors calculated for real and generated molecules. The FCD is arguably the most important metric since it captures the similarity between real and generated molecules' chemical properties. For the generated generic graphs, we evaluate the MMD distance between the test graphs and the generated graphs in terms of four type of graph features, including the statistics of degrees (Deg.), clustering coefficients (Clus.), and orbit counts (Orbit) from You et al. (2018), and the neural-based graph embedding (NN) proposed by Thompson et al. (2022). We report the metrics of the molecule and generic graph generation using 10000 and 128 samples, respectively.

### 5.2 Quantitative Analysis

Here we quantitatively estimate the quality of the generated molecules and generic graphs and report them in Table 1 and Table 2, respectively. We also analyze the sampling speed of NVDiff using the Ego dataset.
**Molecule generation.** NVDiff outperforms all baselines in terms of the FCD metric, thus superior in generating valid molecules that share highly similar characteristics with real molecules. However, on ZINC250K, we observe that NVDiff exhibits relatively lower validity. We argue that filtering invalid molecules is easy as it can be done by checking the valency rule while selecting realistic ones is difficult. We also observe that the sampling speed of NVDiff is much faster than most baselines.
**Generic graph generation.** We show the MMD evaluation on four datasets. Community and Ego contain massive nodes

Table 2: Performance on generic graph generation. Numbers in bold indicate that the methods are the best, numbers underlined are the second best. "OOM" denotes out of memory.

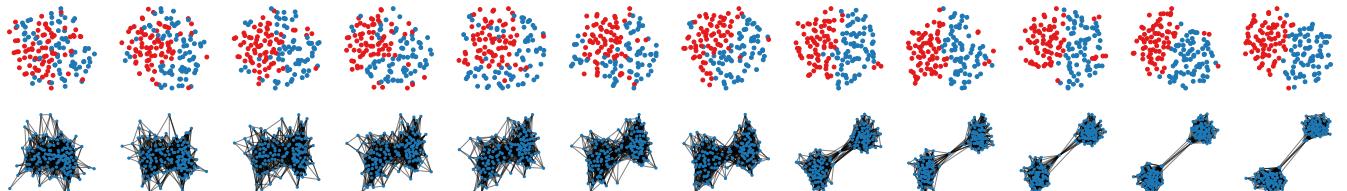| | Community-small | | | | Ego-small | | | |
|---|---|---|---|---|---|---|---|---|
| | Deg.↓ | Clus.↓ | Orbit↓ | NN↓ | Deg.↓ | Clus.↓ | Orbit↓ | NN↓ |
| GraphRNN | 0.080 | 0.120 | 0.040 | 0.997 | 0.090 | 0.220 | <u>0.003</u> | **1.094** |
| GRAN | 0.070 | <u>0.045</u> | 0.021 | <u>0.701</u> | 0.020 | 0.126 | 0.010 | **1.095** |
| GraphCNF | **0.021** | 0.141 | 0.044 | 0.876 | <u>0.011</u> | **0.011** | **0.001** | **1.094** |
| EDP-GNN | 0.053 | 0.144 | 0.026 | 0.789 | 0.052 | 0.093 | 0.007 | **1.095** |
| GDSS | <u>0.045</u> | 0.086 | **0.007** | 0.858 | 0.021 | <u>0.024</u> | 0.007 | <u>1.097</u> |
| NVDiff | **0.021** | **0.035** | <u>0.018</u> | **0.688** | **0.005** | 0.045 | **0.001** | **1.092** |
| | Community | | | | Ego | | | |
| | Deg.↓ | Clus.↓ | Orbit↓ | NN↓ | Deg.↓ | Clus.↓ | Orbit↓ | NN↓ |
| GraphRNN | **0.014** | **0.002** | <u>0.039</u> | 1.452 | 0.077 | 0.316 | <u>0.030</u> | 0.192 |
| GRAN | <u>0.085</u> | <u>0.068</u> | 0.042 | 1.406 | <u>0.064</u> | <u>0.279</u> | 0.039 | 0.196 |
| GraphCNF | 0.707 | 1.305 | 0.194 | 1.356 | OOM | OOM | OOM | OOM |
| EDP-GNN | 0.137 | 1.104 | 0.051 | <u>1.146</u> | 0.069 | 0.447 | 0.055 | **0.189** |
| GDSS | 0.099 | 0.327 | 0.162 | 1.203 | 0.134 | 0.372 | 0.139 | 0.295 |
| NVDiff | 0.093 | 0.086 | **0.021** | **0.982** | **0.045** | **0.091** | **0.004** | 0.195 |



Figure 2: Visualization of the denoising dynamics on Community dataset (evolve from left to right). Top row: t-SNE visualization over the node vectors, with color indicating the community membership. Bottom row: graphs generated from the corresponding node vectors via the decoder.
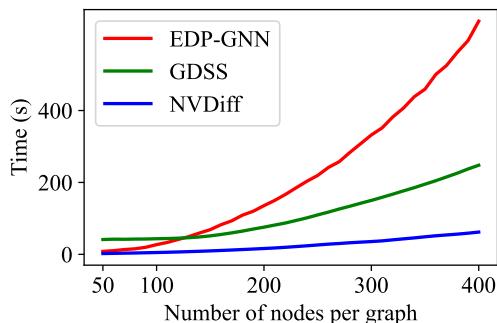


Figure 3: Sampling speed

and thus are difficult to model by several baseline methods. We only report performance for baselines that converged in a reasonable time for these two datasets without consuming extensive computational resources. We observe that NVDiff outperforms all the baselines on most metrics. In GraphRNN, data are fit using BFS orderings, which are favored by datasets with community structure–the model naturally learns to generate communities sequentially. We argue that even without introducing such inductive bias, NVDiff still can capture the community pattern. In comparison, EDP-GNN and GraphCNF fail to capture the pattern of Community and Ego or exhibit poor performance. We hypothesize that

the node-edge relationship becomes more complex as the graph size increases. So modeling nodes and edges in the diffusion prior is not a proper choice. Instead, the diffusion prior in NVDiff only needs to learn the node-to-node relationship, leaving the task of composing edges to the decoder.

**Sampling speed.** Here we compared against other score-based graph generative models in terms of the sampling speed. We choose Ego for runtime analysis. Specifically, we manually fix the number of nodes the model will generate. For each setting, we sample 128 graphs and record the wall clock time. Figure 3 shows the running time concerning the number of nodes of the generated graphs. NVDiff generates much faster than EDP-GNN and GDSS, especially over larger graphs.

### 5.3 Qualitative Analysis

We examine the effectiveness of NVDiff from three aspects: the denoising dynamics of node vectors, the quality of the generated molecules, and the graph statistics encoded in the contextual vector.

**Denoising dynamics of node vectors.** Figure 2 shows how NVDiff generates node vectors and forms a community graph progressively. We record the node vectors along the denoising step. For each time step, we use t-SNE (Van der Maaten and Hinton 2008) to compress the node vectors for visualization. We color each node with its community membership obtained from the spectrum clustering. For each intermediate node vectors, we generate the graph using the learned decoder.

Figure 4: Molecule visualization with Tanimoto similarity. The left-most column visualizes the reference molecules. The top 2 molecules ranked by the similarity scores for each model are visualized. Extended results are in the Appendix.
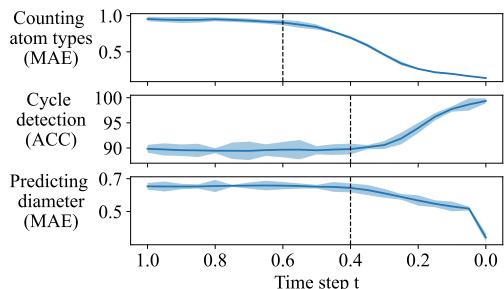


Figure 5: Downstream tasks with contextual vector

As the denoising step increases, node vectors from different communities become separable.

**Sample quality.** We inspect the drug-likeness of the generated samples by measuring the Tanimoto similarity of the Morgan fingerprint (Bajusz, Rácz, and Héberger 2015). Specifically, we calculate the pair-wise Tanimoto similarity between each generated sample and a real molecule randomly chosen from the dataset. Given a reference molecule, we retrieve the top 2 generated molecules ranked by the similarity score and visualize them in Figure 4. We analyze samples generated from NVDiff, GraphDF, Moflow, and GraphCNF. Compared to other baselines, NVDiff can generate molecules that share a large substructure with the actual molecules, illustrating a superior modeling capability.

**Graph statistics in contextual vector.** We investigate whether the introduced contextual vector contains meaningful information about the graph, which can be used for supervising the diffusion process. Here we focus on QM9. We hypothesize that the contextual vector contains both the node and structural information. To validate this, we design three downstream tasks: atom type counting, cycle detection, and diameter prediction. The first task focuses on the global node statistics, and the last two tasks focus on the local and global structural contexts, respectively. We report accuracy for the cycle detection task and Mean Absolute Error (MAE) for the others. Figure 5 demonstrates the downstream task performance of the contextual vectors extracted from different time steps. We can observe that all task performance improves incrementally as time step decreases (i.e., node vectors get less noisy). Moreover, the predictions are precise when $t$ is close to 0, indicating that the contextual vector encodes significant contexts locally and globally about the graph. Interestingly, we find that the performance of tasks related to node context starts to converge earlier than those related to structural context. One reasonable explanation is that the node context is formed earlier than the structural context, as the latter usually depends on the former. We highlight the time step when the task starts to converge in the figure with a dash line.

## 6 Conclusions and limitations

This paper introduces NVDiff that generates novel and realistic graphs by taking the VGAE structure and uses SGM as its prior for latent node vectors. NVDiff guarantees exchangeability without losing expressiveness by implementing SGM with attention-based networks. When pairing with a properly designed GNN decoder, SGM can sample complex graphs in one shot and enjoys a fast sampling speed.

This paper has the following limitations: (1) NVDiff optimizes the variational lower bound of the likelihood, which corresponds to the KL divergence between the empirical distribution and the parameterized distribution. Optimization under other objectives is not discussed. (2) Current NVDiff does not support conditional graph generation, which can be used for optimizing molecular property.

# References

Aldous, D. J. 1981. Representations for partially exchangeable arrays of random variables. *Journal of Multivariate Analysis*, 11(4): 581–598.

Anderson, B. D. 1982. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3): 313–326.

Bacciu, D.; Micheli, A.; and Podda, M. 2020. Edge-based sequential graph generation with recurrent neural networks. *Neurocomputing*, 416: 177–189.

Bajusz, D.; Rácz, A.; and Héberger, K. 2015. Why is Tanimoto index an appropriate choice for fingerprint-based similarity calculations? *Journal of cheminformatics*, 7(1): 1–13.

Bengio, E.; Jain, M.; Korablyov, M.; Precup, D.; and Bengio, Y. 2021. Flow network based generative models for non-iterative diverse candidate generation. *Advances in Neural Information Processing Systems*, 34: 27381–27394.

Bradshaw, J.; Paige, B.; Kusner, M. J.; Segler, M.; and Hernández-Lobato, J. M. 2019. A model to search for synthesizable molecules. *Advances in Neural Information Processing Systems*, 32.

Chen, N.; Zhang, Y.; Zen, H.; Weiss, R. J.; Norouzi, M.; and Chan, W. 2020. WaveGrad: Estimating gradients for waveform generation. *arXiv preprint arXiv:2009.00713*.

Chen, X.; Han, X.; Hu, J.; Ruiz, F. J.; and Liu, L. 2021. Order Matters: Probabilistic Modeling of Node Sequence for Graph Generation. *arXiv preprint arXiv:2106.06189*.

Costa, F.; and De Grave, K. 2010. Fast neighborhood subgraph pairwise distance kernel. In *ICML*.

Dai, B.; and Wipf, D. 2019. Diagnosing and enhancing VAE models. *arXiv preprint arXiv:1903.05789*.

Dai, H.; Nazi, A.; Li, Y.; Dai, B.; and Schuurmans, D. 2020. Scalable deep generative modeling for sparse graphs. In *International Conference on Machine Learning*, 2302–2312. PMLR.

Garcia Satorras, V.; Hoogeboom, E.; Fuchs, F.; Posner, I.; and Welling, M. 2021. E (n) Equivariant Normalizing Flows. *Advances in Neural Information Processing Systems*, 34.

Goyal, N.; Jain, H. V.; and Ranu, S. 2020. GraphGen: a scalable approach to domain-agnostic labeled graph generation. In *Proceedings of The Web Conference 2020*, 1253–1263.

Gretton, A.; Borgwardt, K. M.; Rasch, M. J.; Schölkopf, B.; and Smola, A. 2012. A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1): 723–773.

Grover, A.; Zweig, A.; and Ermon, S. 2019. Graphite: Iterative generative modeling of graphs. In *International conference on machine learning*, 2434–2444. PMLR.

Guo, X.; Zhao, L.; Qin, Z.; Wu, L.; Shehu, A.; and Ye, Y. 2020. Node-edge co-disentangled representation learning for attributed graph generation. In *International Conference on Knowledge Discovery and Data Mining (SIGKDD)*.

Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33: 6840–6851.

Hoover, D. 1979. Relations on probability spaces and arrays of random variables. *Preprint*.

Irwin, J. J.; Sterling, T.; Mysinger, M. M.; Bolstad, E. S.; and Coleman, R. G. 2012. ZINC: a free tool to discover chemistry for biology. *Journal of chemical information and modeling*, 52(7): 1757–1768.

Jin, W.; Barzilay, R.; and Jaakkola, T. 2018. Junction tree variational autoencoder for molecular graph generation. In *International conference on machine learning*, 2323–2332. PMLR.

Jin, W.; Barzilay, R.; and Jaakkola, T. 2020. Hierarchical generation of molecular graphs using structural motifs. In *International Conference on Machine Learning*, 4839–4848. PMLR.

Jo, J.; Lee, S.; and Hwang, S. J. 2022. Score-based Generative Modeling of Graphs via the System of Stochastic Differential Equations. *arXiv preprint arXiv:2202.02514*.

Kajino, H. 2019. Molecular hypergraph grammar with its application to molecular optimization. In *International Conference on Machine Learning*, 3183–3191. PMLR.

Kingma, D. P.; and Welling, M. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.

Kipf, T. N.; and Welling, M. 2016. Variational graph autoencoders. *arXiv preprint arXiv:1611.07308*.

Klushyn, A.; Chen, N.; Kurle, R.; Cseke, B.; and van der Smagt, P. 2019. Learning hierarchical priors in vaes. *Advances in neural information processing systems*, 32.

Li, J.; Yu, J.; Li, J.; Zhang, H.; Zhao, K.; Rong, Y.; Cheng, H.; and Huang, J. 2020. Dirichlet graph variational autoencoder. *Advances in Neural Information Processing Systems*, 33: 5274–5283.

Li, Y.; Vinyals, O.; Dyer, C.; Pascanu, R.; and Battaglia, P. 2018. Learning deep generative models of graphs. *arXiv preprint arXiv:1803.03324*.

Li, Y.; Zhang, L.; and Liu, Z. 2018. Multi-objective de novo drug design with conditional graph generative model. *Journal of cheminformatics*, 10(1): 1–24.

Liao, R.; Li, Y.; Song, Y.; Wang, S.; Hamilton, W.; Duvenaud, D. K.; Urtasun, R.; and Zemel, R. 2019. Efficient graph generation with graph recurrent attention networks. *Advances in Neural Information Processing Systems*, 32.

Lim, J.; Hwang, S.-Y.; Moon, S.; Kim, S.; and Kim, W. Y. 2020. Scaffold-based molecular design with a graph generative model. *Chemical science*, 11(4): 1153–1164.

Lippe, P.; and Gavves, E. 2020. Categorical normalizing flows via continuous transformations. *arXiv preprint arXiv:2006.09790*.

Liu, J.; Kumar, A.; Ba, J.; Kiros, J.; and Swersky, K. 2019. Graph normalizing flows. *Advances in Neural Information Processing Systems*, 32.

Liu, Q.; Allamanis, M.; Brockschmidt, M.; and Gaunt, A. 2018. Constrained graph variational autoencoders for molecule design. *Advances in neural information processing systems*, 31.

Luo, Y.; Yan, K.; and Ji, S. 2021. GraphDF: A discrete flow model for molecular graph generation. In *International Conference on Machine Learning*, 7192–7203. PMLR.

Madhawa, K.; Ishiguro, K.; Nakago, K.; and Abe, M. 2019. Graphnvp: An invertible flow model for generating molecular graphs. *arXiv preprint arXiv:1905.11600*.

Mehta, N.; Duke, L. C.; and Rai, P. 2019. Stochastic block-models meet graph neural networks. In *International Conference on Machine Learning*, 4466–4474. PMLR.

Mittal, G.; Engel, J.; Hawthorne, C.; and Simon, I. 2021. Symbolic music generation with diffusion models. *arXiv preprint arXiv:2103.16091*.

Niu, C.; Song, Y.; Song, J.; Zhao, S.; Grover, A.; and Ermon, S. 2020. Permutation invariant graph generation via score-based generative modeling. In *International Conference on Artificial Intelligence and Statistics*, 4474–4484. PMLR.

Preuer, K.; Renz, P.; Unterthiner, T.; Hochreiter, S.; and Klambauer, G. 2018. Fréchet ChemNet distance: a metric for generative models for molecules in drug discovery. *Journal of chemical information and modeling*, 58(9): 1736–1741.

Ramakrishnan, R.; Dral, P. O.; Rupp, M.; and Von Lilienfeld, O. A. 2014. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific data*, 1(1): 1–7.

Rezende, D.; and Mohamed, S. 2015. Variational inference with normalizing flows. In *International conference on machine learning*, 1530–1538. PMLR.

Samanta, B.; De, A.; Jana, G.; Gómez, V.; Chattaraj, P. K.; Ganguly, N.; and Gomez-Rodriguez, M. 2020. Nevae: A deep generative model for molecular graphs. *Journal of machine learning research. 2020 Apr; 21 (114): 1-33*.

Satorras, V. G.; Hoogeboom, E.; and Welling, M. 2021. E (n) equivariant graph neural networks. In *International Conference on Machine Learning*, 9323–9332. PMLR.

Shi, C.; Xu, M.; Zhu, Z.; Zhang, W.; Zhang, M.; and Tang, J. 2020. Graphaf: a flow-based autoregressive model for molecular graph generation. *arXiv preprint arXiv:2001.09382*.

Simonovsky, M.; and Komodakis, N. 2018. Graphvae: Towards generation of small graphs using variational autoencoders. In *International conference on artificial neural networks*, 412–422. Springer.

Sinha, A.; Song, J.; Meng, C.; and Ermon, S. 2021. D2C: Diffusion-Decoding Models for Few-Shot Conditional Generation. *Advances in Neural Information Processing Systems*, 34.

Song, Y.; and Ermon, S. 2019. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32.

Song, Y.; Sohl-Dickstein, J.; Kingma, D. P.; Kumar, A.; Ermon, S.; and Poole, B. 2020. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*.

Thompson, R.; Knyazev, B.; Ghalebi, E.; Kim, J.; and Taylor, G. W. 2022. On Evaluation Metrics for Graph Generative Models. *arXiv preprint arXiv:2201.09871*.

Tran, C.; Shin, W.-Y.; Spitz, A.; and Gertz, M. 2020. DeepNC: Deep generative network completion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

Vahdat, A.; Kreis, K.; and Kautz, J. 2021. Score-based generative modeling in latent space. *Advances in Neural Information Processing Systems*, 34.

Van der Maaten, L.; and Hinton, G. 2008. Visualizing data using t-SNE. *Journal of machine learning research*, 9(11).

Watts, D. J.; and Strogatz, S. H. 1998. Collective dynamics of 'small-world' networks. *nature*, 393(6684): 440–442.

Wehenkel, A.; and Louppe, G. 2021. Diffusion priors in variational autoencoders. *arXiv preprint arXiv:2106.15671*.

Xu, M.; Luo, S.; Bengio, Y.; Peng, J.; and Tang, J. 2021. Learning neural generative dynamics for molecular conformation generation. *arXiv preprint arXiv:2102.10240*.

You, J.; Ying, R.; Ren, X.; Hamilton, W.; and Leskovec, J. 2018. Graphrnn: Generating realistic graphs with deep autoregressive models. In *International conference on machine learning*, 5708–5717. PMLR.

Zang, C.; and Wang, F. 2020. MoFlow: an invertible flow model for generating molecular graphs. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 617–626.

## Implementation Details

### Encoder

Here we demonstrate the architecture of the encoder $\text{gnn}_\phi(\cdot, \cdot)$. We treat non-edge as an extra edge type and use a GNN to output the encoding distribution. First, we create the input node representation $\mathbf{M}^0$ for the encoder by injecting Gaussian Noise

$$\mathbf{M}^0 = [\mathbf{X}; \boldsymbol{\epsilon}^{\text{in}}], \boldsymbol{\epsilon}^{\text{in}} \sim \mathcal{N}(\boldsymbol{\epsilon}^{\text{in}}; \mathbf{0}, \mathbf{I}), \tag{16}$$

Denote $\mathbf{m}_i^\ell$ to be the $i$-th row of $\mathbf{M}^\ell$, the $\ell$-th layer of the encoder updates the node representation as follow:

$$\mathbf{m}_i^\ell = \mathbf{m}_i^{\ell-1} + \text{mlp}^v([\mathbf{m}_i^\ell; \sum_{j \neq i}^N \mathbf{m}_{i,j}^\ell]),$$
$$\mathbf{m}_{i,j}^\ell = \text{mlp}^e([\mathbf{a}_{i,j}; \mathbf{W}\mathbf{m}_i^{\ell-1} + \mathbf{W}\mathbf{m}_j^{\ell-1}]). \tag{17}$$

Let $L'$ be the number of layer of the $\text{gnn}_\phi(\cdot, \cdot)$. We denote $\mathbf{M} := \mathbf{M}^{L'}$ being the mean of the encoding distribution.

### NVDiff-E

To demonstrate the effectiveness of the introduced contextual vector, we propose a variant of the score network that diffuses both nodes and edges, which we call NVDiff-E. Compared to NVDiff, NVDiff-E exhibits a slower sampling speed and is incapable of modeling larger graphs. In molecule generation, where graphs are relatively small, NVDiff-E performs better than most methods, including NVDiff. Since NVDiff-E can operate directly on data space $(\mathbf{A}, \mathbf{X})$, we only show how its score network is built.

The score network for NVDiff-E $\epsilon_\theta^{\text{AX}}(\mathbf{A}^t, \mathbf{X}^t, t)$ takes in the noisy inputs $(\mathbf{A}^t, \mathbf{X}^t)$ and time step $t \in [0, 1]$, and predict the scores for the edge and node features, respectively. We first inject time step information into the node and edge inputs:

$$\mathbf{H}_0^e = \text{mlp}^{\text{in},e}([\mathbf{A}^t; \mathbf{1}\mathbf{t}^\top]),$$
$$\mathbf{H}_0^v = \text{mlp}^{\text{in},v}([\mathbf{X}^t; \mathbf{1}\mathbf{t}^\top]), \mathbf{t} = \text{pe}(t). \tag{18}$$

Here we define the edge-node attention layer (ENA). We stack $L$ ENA layers to obtain the final representations $(\mathbf{H}_L^e, \mathbf{H}_L^v)$. We defer the details of the ENA module, and the following computations are repeated $L$ times:

$$(\mathbf{g}_\ell, \mathbf{H}_\ell^v, \mathbf{H}_\ell^e) = \text{ENA}(\mathbf{g}_{\ell-1}, \mathbf{H}_{\ell-1}^v, \mathbf{H}_{\ell-1}^e), \ell = 1, \ldots, L. \tag{19}$$

Here $\mathbf{g}_\ell$ is the contextual embedding. We predict the scores from the final representations.

$$\hat{\boldsymbol{\epsilon}}^e = \text{mlp}^{\text{out},e}(\mathbf{H}_L^e), \hat{\boldsymbol{\epsilon}}^v = \text{mlp}^{\text{out},v}(\mathbf{H}_L^v) \tag{20}$$

**Edge-node attention.** Here we define ENA module. The ENA module is used to obtain $(\mathbf{g}_\ell, \mathbf{H}_\ell^v, \mathbf{H}_\ell^e)$ from $(\mathbf{g}_{\ell-1}, \mathbf{H}_{\ell-1}^v, \mathbf{H}_{\ell-1}^e)$:

$$(\mathbf{g}_\ell', \mathbf{H}_\ell^{v'}) = \text{mha}(\mathbf{g}_{\ell-1}, \mathbf{H}_{\ell-1}^v), \mathbf{g}_\ell = \mathbf{g}_{\ell-1} + \mathbf{g}_\ell',$$
$$\mathbf{h}_{i,\ell}^v = \mathbf{h}_{i,\ell}^{v}{}' + \sum_{j:j \neq i} \mathbf{W}^v \mathbf{h}_{i,j,\ell-1}^e,$$
$$\mathbf{h}_{i,j,l}^e = \text{mlp}^e([\mathbf{W}^e \mathbf{h}_{i,\ell}^v + \mathbf{W}^e \mathbf{h}_{j,\ell}^v; \mathbf{h}_{i,j,\ell-1}^e]). \tag{21}$$

Here $\mathbf{h}_{i,\ell}^v$ and $\mathbf{h}_{i,\ell}^{v}{}'$ are the $i$-th row in $\mathbf{H}_\ell^v$ and $\mathbf{H}_\ell^{v'}$, respectively. And $\mathbf{h}_{i,j,\ell}^e$ is the $(i, j)$-th vector slice in $\mathbf{H}_\ell^e$.

To draw a sample from NVDiff-E, we first sample $\mathbf{A}^1$ and $\mathbf{X}^1$ from $\mathcal{N}(\mathbf{A}^1; \mathbf{0}, \mathbf{1})$ and $\mathcal{N}(\mathbf{X}^1; \mathbf{0}, \mathbf{1})$, respectively, and run the reverse-time SDE to obtain $(\mathbf{A}^0, \mathbf{X}^0)$. Then we obtain $(\mathbf{A}, \mathbf{X})$ by discretizing $(\mathbf{A}^0, \mathbf{X}^0)$:

$$\mathbf{A} = \text{one\_hot}(\text{argmax}(\mathbf{A}^0)),$$
$$\mathbf{X} = \text{one\_hot}(\text{argmax}(\mathbf{X}^0)). \tag{22}$$

### Cross Entropy

Directly training the KL term in Eq. (2) is impossible as it involves the marginal score $\nabla_{\mathbf{Z}^t} \log q_t(\mathbf{Z}^t)$. While the entropy term is fixed to a constant, we show that cross entropy term can be estimated using the following theorem:

**Theorem 1.** *(Vahdat, Kreis, and Kautz 2021) Given two distributions $q(\mathbf{Z}^0|\mathbf{A}, \mathbf{X})$ and $p(\mathbf{Z}^0)$, defined in the continuous space $\mathbb{R}^{N \times d}$, denote the marginal distributions of diffused samples under the SDE in Eqn. (1) at time $t$ with $q(\mathbf{Z}^t|\mathbf{A}, \mathbf{X})$ and $p(\mathbf{Z}^t)$. Assuming mild smoothness conditions on $\log q(\mathbf{Z}^t|\mathbf{A}, \mathbf{X})$ and $\log p(\mathbf{Z}^t)$, the cross entropy is:*

$$\mathbb{CE}(q(\mathbf{Z}^0|\mathbf{A}, \mathbf{X})||p(\mathbf{Z}^0))$$
$$= \mathbb{E}_{t,\boldsymbol{\epsilon},q_\phi(\mathbf{Z}^0,\mathbf{Z}^t|\mathbf{A},\mathbf{X})} \left[ \frac{g(t)^2}{2} ||\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{Z}^t, t)||_2^2 \right] + \frac{C}{2}, \tag{23}$$

*where $t \sim U[0, 1]$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I})$, $C = Nd \log 2\pi e \sigma_0^2$, and $q_\phi(\mathbf{Z}^0, \mathbf{Z}^t|\mathbf{A}, \mathbf{X}) = q(\mathbf{Z}^t|\mathbf{Z}^0)q_\phi(\mathbf{Z}^0|\mathbf{A}, \mathbf{X})$.*

Here we focus on variance preserving SDE (VPSDE) (Ho, Jain, and Abbeel 2020; Song et al. 2020), and denote $q(\mathbf{Z}^t|\mathbf{Z}^0) = \mathcal{N}(\mathbf{Z}^t; \boldsymbol{\mu}_t(\mathbf{Z}^0), \sigma_t\mathbf{I})$, with $\sigma_0 = 0$ at time step $t = 0$. Specifically, the Gaussian parameters in the transition kernel is computed as follow:

$$\boldsymbol{\mu}_t(\mathbf{Z}^0) = e^{\frac{1}{2}\int_0^t \beta(s)\text{d}s}\mathbf{Z}^0,$$
$$\sigma_t = 1 - (1 - \sigma_0^2)e^{\frac{1}{2}\int_0^t \beta(s)\text{d}s},$$
$$\beta(t) = \beta_0 + (\beta_1 - \beta_0)t, \tag{24}$$

where $\beta_0$ and $\beta_1$ are hyperparameters denoting the variances at the start and the end of the diffusion process, respectively. For the drift and diffusion coefficients in Eqn. (1), we have $f(t) = -\beta(t)/2$ and $g(t) = \sqrt{\beta(t)}$.

Moreover, as we use VPSDE and have $\sigma_0 = 0$ at time step $t = 0$, we are not able to compute the integral of the cross entropy over the full time interval $[0, 1]$. The constant term $Nd/2 \log 2\pi e \sigma_0^2$ Following Vahdat, Kreis, and Kautz (2021), we approximate the cross entropy term by limiting $t \sim U[\epsilon^t, 1]$, where $\epsilon^t$ is another hyperparameter, and yield:

$$\mathbb{CE}(q(\mathbf{Z}^0|\mathbf{A}, \mathbf{X})||p(\mathbf{Z}^0))$$
$$\approx \mathbb{E}_{t,\boldsymbol{\epsilon},q_\phi(\mathbf{Z}^0,\mathbf{Z}^t|\mathbf{A},\mathbf{X})} \left[ \frac{g(t)^2}{2} ||\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{Z}^t, t)||_2^2 \right] + \frac{C_{\epsilon^t}}{2}, \tag{25}$$

where $t \sim U[\epsilon^t, 1]$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I})$, $C_{\epsilon^t} = Nd \log 2\pi e \sigma_{\epsilon^t}^2$, and $q_\phi(\mathbf{Z}^0, \mathbf{Z}^t|\mathbf{A}, \mathbf{X}) = q(\mathbf{Z}^t|\mathbf{Z}^0)q_\phi(\mathbf{Z}^0|\mathbf{A}, \mathbf{X})$. Since the

Table 3: Hyperparameters of NVDiff used in the experiments.

| | QM9 | ZINC250K | Community | Ego | Community-small | Ego-small |
|---|---|---|---|---|---|---|
| **Encoder** | | | | | | |
| # layers | 3 | 5 | 3 | 3 | 3 | 3 |
| # hidden dims | 64 | 64 | 64 | 64 | 32 | 32 |
| # latent dims | 16 | 32 | 8 | 4 | 4 | 4 |
| # noise dims | 8 | 8 | 8 | 8 | 8 | 8 |
| fixed var. | 0.01 | 0.0025 | 0.01 | 0.01 | 0.01 | 0.01 |
| **Decoder** | | | | | | |
| # layers | 3 | 3 | 1 | 1 | 1 | 1 |
| # hidden dims | 64 | 64 | 64 | 64 | 32 | 32 |
| # finetune epochs | 200 | 200 | 0 | 0 | 0 | 0 |
| noise var. in finetuning | 0.0001 | 0.025 | 0 | 0 | 0 | 0 |
| **SGM prior** | | | | | | |
| $\beta_0$ | | | 0.1 | | | |
| $\beta_1$ | | | 20 | | | |
| $\epsilon^t$ | | | 0.01 | | | |
| time embedding | | | positional | | | |
| # time emb. dims | | | 16 | | | |
| SDE type | | | variance preserving SDE | | | |
| # layers | 3 | 5 | 3 | 3 | 3 | 3 |
| # hidden dims | 64 | 64 | 32 | 32 | 16 | 16 |
| # attention heads | 4 | 4 | 2 | 4 | 2 | 4 |
| **Training** | | | | | | |
| optimizer | | | Adam | | | |
| learning rate (VAE) | 1e-4 | 1e-4 | 1e-4 | 1e-4 | 1e-3 | 1e-3 |
| weight decay (VAE) | | | 1e-4 | | | |
| learning rate (SGM) | 1e-4 | 1e-4 | 1e-4 | 1e-4 | 1e-3 | 1e-3 |
| weight decay (SGM) | | | 1e-4 | | | |
| KL annealing to | 0.7 | 0.7 | 1.0 | 1.0 | 1.0 | 1.0 |
| # epochs | 1000 | 2000 | 15000 | 15000 | 4000 | 4000 |
| batch size | 256 | 256 | 4 | 4 | 8 | 8 |
| total training time (h) | 72 | 120 | 72 | 72 | 8 | 8 |
| **Evaluation** | | | | | | |
| # used samples | 10000 | 10000 | 128 | 128 | 128 | 128 |
| ODE solver error tolerance | 1e-4 | 1e-4 | 1e-4 | 1e-4 | 1e-5 | 1e-5 |

SGM prior never learns to recover $\mathbf{Z}^0$ exactly, the stage for finetuning the decoder by injecting extra noise becomes necessary when the sample quality heavily relies on the reconstruction accuracy. Empirically, we find that the finetune stage improves the validity metric by 15% on the ZINC250K dataset.

## Exchangeability of NVDiff

We guarantee the encoder and the decoder are permutation-equivariant by utilizing GNNs. And since the SGM prior is designed to be an exchangeable distribution, NVDiff is guaranteed to be an exchangeable model.

## Dataset Statistics

Table 4 shows the details of the used datasets. We use 90% of each dataset for training and test the rest for molecular datasets. For generic graph datasets, we use 80% of the dataset for training.

## Experimental Details

All hyperparameters of NVDiff on different datasets are provided in Table 3.

## Training

We mostly follow the previous works to optimize NVDiff. We pretrain the VAE with a fixed prior before warming up the encoder and decoder before replacing the SGM prior. During the pretraining, KL annealing is also applied. We adopt the importance sampling scheme proposed by (Vahdat, Kreis, and Kautz 2021) to sample time step t to reduce the gradient variance, thus enabling faster convergence. Following Sinha et al. (2021), we normalize the latent node vectors obtained from the encoder before feeding it into the SGM prior. We did not observe improvement by utilizing an exponential moving average (EMA) for SGM prior and VAE parameters during test time. We also apply gradient clipping during the optimization. For the edge distribution in the decoder, we optimize the edge-existence distribution $p(\mathbf{a}_{i,j}|\mathbf{r}_{i,j})$ and the edge-type

Table 4: Dataset statistics

| | Molecular dataset | | Generic graph dataset | | | |
| | QM9 | ZINC250K | Community | Ego | Community-small | Ego-small |
|---|---|---|---|---|---|---|
| # graphs | 130828 | 249456 | 500 | 753 | 500 | 500 |
| # nodes(min,max) | (2, 9) | (6, 38) | (60, 160) | (50, 399) | (12, 20) | (4, 18) |
| # edges(min,max) | (1, 14) | (5, 45) | (240, 1995) | (57, 1071) | (23, 82) | (3, 63) |
| # node types | 4 | 9 | 1 | 1 | 1 | 1 |
| # edge types | 4 | 4 | 1 | 1 | 1 | 1 |

Table 5: Further experiment results on molecule generation. Results of baselines marked with "*" are taken from previous papers. We used paired t-test to compare the results; the numbers in bold indicate that the method is better at the 5% significance level; the number underlined are the second best.

| | QM9 | | | |
| | Validity↑ | Uniqueness↑ | NSPDK↓ | FCD↓ |
|---|---|---|---|---|
| GraphDF* | 82.67 | 97.62 | 0.063±0.001 | 10.816±0.020 |
| Moflow* | 91.36±1.23 | 94.72±0.77 | 0.017±0.003 | 4.467±0.595 |
| GraphCNF | 95.00±1.15 | 93.68±0.96 | **0.002±0.000** | 1.629±0.326 |
| EDP-GNN* | 47.52±3.60 | **99.25±0.05** | 0.005±0.001 | 2.680±0.221 |
| GDSS* | 95.72±1.94 | 98.46±0.61 | 0.003±0.000 | 2.900±0.282 |
| NVDiff | 95.79±1.32 | 96.53±1.16 | **0.002±0.000** | 1.131±0.213 |
| NVDiff(10%) | 90.29±1.12 | 94.88±0.96 | 0.003±0.000 | 1.758±0.390 |
| NVDiff-E | **98.49±0.81** | 98.32±0.55 | **0.002±0.000** | **0.350±0.122** |
| | ZINC250K | | | |
| | Validity↑ | Uniqueness↑ | NSPDK↓ | FCD↓ |
| GraphDF* | 89.03 | 99.16 | 0.176±0.001 | 34.202±0.160 |
| Moflow* | 63.11±5.17 | **100±0.00** | 0.046±0.002 | 20.931±0.184 |
| GraphCNF | **95.26±1.85** | 99.96±0.03 | 0.089±0.002 | 14.786±1.326 |
| EDP-GNN* | 82.97±2.73 | 99.79±0.08 | 0.049±0.006 | 16.737±1.300 |
| GDSS* | **97.01±0.77** | 99.64±0.13 | **0.019±0.001** | 14.656±0.680 |
| NVDiff | 85.63±0.89 | **99.99±0.01** | 0.066±0.006 | **4.019±0.01** |
| NVDiff-E | **96.23±1.21** | 99.94±0.03 | 0.061±0.001 | 6.681±0.01 |

distribution $p(\mathbf{a}_{i,j}|\mathbf{a}_{i,j} \neq \mathbf{0}, \mathbf{r}_{i,j})$ at the same time. Specifically, we fit the edge-existence distribution with the binary adjacency matrix and the edge-type distribution for all edge entries weighted by the ground-truth edge-existence indicator. During finetuning, we fix the parameters of the encoder and only train the decoder for more epochs by optimizing $\mathcal{L}_{\mathrm{Reg}}$.

## Evaluation

We report all results over 5 runs with different random seeds for each dataset.

**Molecule generation.** We precompute the prior distribution of graph size over the training set. The graph size distribution is then used to draw the number of nodes each sample should contain. We report Uniqueness, NSPDK MMD, and FCD metrics over the valid samples. We do not employ any validity correction scheme during training or sampling (Shi et al. 2020; Luo, Yan, and Ji 2021).

**Generic graph generation.** We compute MMD between the 128 generated graphs and graphs in the test set. We mainly follow the configurations in You et al. (2018); Thompson et al. (2022) when evaluating MMD. We extract the largest sub-component of the generated graphs when evaluating Ego and Ego-small.

## Configuration of Downstream tasks for Contextual Vector

We use QM9 for contextual vector analysis. Here we specify the setups for each downstream tasks. We record the contextual vectors over all training graphs for t = [0, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95, 1.0]. For each intermediate time step t, We train and evaluate on the collected contextual vectors. We use 90% of the data for training and 10% of the data for testing. For all tasks, we use a 3-layer MLP.

**Counting atom types.** QM9 contains four heavy atoms (O, N, C, F). Here we manually generate the label of the number of atoms of the four types for each molecule. We consider the counting problem as each atom type as an independent task. We report the MAE metric by averaging over four counting tasks.

**Cycle detection and diameter prediction.** We only consider a binary classification task for cycle detection and use an MLP to determine whether a molecule contains a cycle. We ignore the atom type and bond type when computing the existence of the cycle and the diameter of the molecules. We perform a regression task on the diameter prediction.

Table 6: Ablation study of contextual vector on QM9 and ZINC250K.

|  | QM9 | | ZINC250K | |
|---|---|---|---|---|
|  | Validity | FCD | Validity | FCD |
| w/o $g_0$ | 93.18 | 3.082 | 83.11 | 6.100 |
| w/ $g_0$ | 95.79 | 1.131 | 85.63 | 4.019 |

Table 7: Ablation study of consistence regularization on QM9 and ZINC250K.

|  | QM9 | | ZINC250K | |
|---|---|---|---|---|
|  | Validity | FCD | Validity | FCD |
| w/o $g_0$ | 94.24 | 1.021 | 71.35 | 4.295 |
| w/ $g_0$ | 95.79 | 1.131 | 85.63 | 4.019 |

## Computational Resources

We train and evaluate NVDiff on GeForce RTX 3060, A100, or RTX 6000 GPU.

# Additional Results

## Additional Results on Generic Graph Generation

We additionally provide the standard derivation of the all the results for generic graph generation in Table 8.

## Additional Metrics on Molecule Generation

We additionally provide the Uniqueness, and NSPDK metrics as well as all standard derivations of the results in Table 5. We also report the performance of the NVDiff-E.

## Ablation Studies

Here, we present two ablation studies investigating the effectiveness of contextual embedding and consistency regularization term, respectively. We examine the two components using molecular datasets.

**Contextual vector.** We train NVDiff with diffusion prior without the contextual embedding $g_0$ on two molecular datasets and report the Validity and FCD scores. From Table 6, we can observe that using contextual embedding improves the validity and FCD metrics.

**Consistency regularization.** The consistency regularization term improves the robustness of the decoder given noisy samples generated by the diffusion prior. Thus, regularization helps improve the reconstruction precision, which is especially valuable in drug synthesis. We report the effect of the consistency regularization term on the validity metrics over the two molecular datasets. As we can see in Table 7, by introducing the consistency regularization term, the model is able to generate molecules with higher validity.

## Visualization of Generated Molecules

Here we visualize the result for eight reference molecules randomly picked from the dataset. We visualize the top 2 generated samples ranked by the similarity score. Figure 6 and Figure 7 show the visualization results.

Table 8: Further experiment results on generic graph generation. Results of baselines marked with "*" are taken from previous papers. Numbers in bold indicate that the methods are the best, numbers underlined are the second best.

| | Deg.↓ | Clus.↓ | Orbit↓ | NN↓ |
|---|---|---|---|---|
| **Community-small** | | | | |
| GraphRNN* | 0.08 | 0.120 | 0.040 | 0.997±0.010 |
| GRAN | 0.070±0.009 | 0.045±0.004 | 0.021±0.003 | 0.701±0.093 |
| GraphCNF | **0.021±0.001** | 0.141±0.004 | 0.044±0.006 | 0.876±0.138 |
| EDP-GNN* | 0.053 | 0.144 | 0.026 | 0.798±0.159 |
| GDSS* | 0.045±0.028 | 0.086±0.022 | **0.007±0.004** | 0.858±0.166 |
| NVDiff | **0.021±0.001** | **0.035±0.002** | 0.018±0.002 | **0.688±0.144** |
| **Community** | | | | |
| GraphRNN* | **0.014** | **0.002** | 0.039 | 1.452±0.092 |
| GRAN | 0.085±0.007 | 0.068±0.008 | 0.042±0.003 | 1.406±0.191 |
| GraphCNF | 0.707±0.129 | 1.305±0.326 | 0.194±0.062 | 1.356±0.261 |
| EDP-GNN | 0.137±0.042 | 1.104±0.036 | 0.051±0.019 | 1.416±0.133 |
| GDSS | 0.099±0.009 | 0.327±0.023 | 0.162±0.057 | 1.203±0.353 |
| NVDiff | 0.093±0.008 | 0.086±0.006 | **0.021±0.001** | **0.982±0.155** |
| **Ego-small** | | | | |
| GraphRNN* | 0.090 | 0.220 | 0.003 | **1.094±0.003** |
| GRAN | 0.020±0.001 | 0.126±0.021 | 0.010±0.002 | **1.095±0.003** |
| GraphCNF | 0.011±0.001 | **0.011±0.001** | **0.001±0.000** | **1.094±0.002** |
| EDP-GNN* | 0.052 | 0.093 | 0.007 | **1.095±0.002** |
| GDSS* | 0.021±0.008 | 0.024±0.007 | 0.007±0.004 | 1.097±0.001 |
| NVDiff | **0.005±0.001** | 0.045±0.004 | **0.001±0.000** | 1.092±0.002 |
| **Ego** | | | | |
| GraphRNN* | 0.077 | 0.316 | 0.030 | 0.192±0.002 |
| GRAN | 0.064±0.012 | 0.279±0.046 | 0.039±0.009 | 0.196±0.008 |
| GraphCNF | OOM | OOM | OOM | OOM |
| EDP-GNN* | 0.069±0.010 | 0.447±0.065 | 0.055±0.014 | **0.189±0.001** |
| GDSS* | 0.134±0.012 | 0.371±0.042 | 0.148±0.02 | 0.295±0.069 |
| NVDiff | **0.045±0.004** | **0.091±0.011** | **0.004±0.001** | 0.195±0.004 |

Figure 6: Extended result of molecule visualization on QM9.

Figure 7: Extended result of molecule visualization on ZINC250K.