

---

# Diffusion Models for Graphs Benefit From Discrete State Spaces

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Denoising diffusion probabilistic models and score matching models have proven  
2 to be very powerful for generative tasks. While these approaches have also been  
3 applied to the generation of discrete graphs, they have, so far, relied on continuous  
4 Gaussian perturbations. Instead, in this work, we suggest using discrete noise for  
5 the forward Markov process. This ensures that in every intermediate step the graph  
6 remains discrete. Compared to the previous approach, our experimental results on  
7 four datasets and multiple architectures show that using a discrete noising process  
8 results in higher quality generated samples indicated with an average MMDs  
9 reduced by a factor of 1.5. Furthermore, the number of denoising steps is reduced  
10 from 1000 to 32 steps leading to a 30 times faster sampling procedure.

## 11 1 Introduction

12 Score-based [1] and denoising diffusion probabilistic models (DDPMs) [2, 3] have recently achieved  
13 striking results in generative modeling and in particular in image generation. Instead of learning a  
14 complex model that generates samples in a single pass (like a Generative Adversarial Network [4]  
15 (GAN) or a Variational Auto-Encoder [5] (VAE)), a diffusion model is a parameterized Markov  
16 Chain trained to reverse an iterative predefined process that gradually transforms a sample into pure  
17 noise. Although diffusion processes have been proposed for both continuous [6] and discrete [7]  
18 state spaces, their use for graph generation has only focused on Gaussian diffusion processes which  
19 operate in the continuous state space [8, 9].

20 This contribution suggests adapting the denoising procedure to an actual graph distribution and  
21 using discrete noise, leading to a random graph model. We describe this procedure based on the  
22 Discrete DDPM framework proposed by Austin et al. [7], Hooigeboom et al. [10]. Our experiments  
23 show that using discrete noise greatly reduces the number of denoising steps that are needed and  
24 improves the sample quality. We also suggest the use of a simple expressive graph neural network  
25 architecture [11] for denoising, which, while bringing expressivity benefits, contrasts with more  
26 complicated architectures currently used for graph denoising [8].

## 27 2 Related Work

28 Traditionally, graph generation has been studied through the lens of random graph models [12–14].  
29 While this approach is insufficient to model many real-world graph distributions, it is useful to create  
30 synthetic datasets and provides a useful abstraction. In fact, we will use Erdős–Rényi (ER) graphs [12]  
31 to model the prior distribution of our diffusion process.

32 Due to their expressive power, deep generative models have achieved better results in modeling  
33 complex graph distributions. The most successful graph generative models can be divided into two

34 camps: a) auto-regressive graph generative models, which generate the graph sequentially node-by-  
 35 node [15, 16], and b) one-shot generative models which generate the whole graph in a single forward  
 36 pass [17–20, 8, 9, 21]. While auto-regressive models can generate graphs with hundreds or even  
 37 thousands of nodes, they can suffer from mode collapse [20, 21]. One-shot graph generative models  
 38 are more resilient to mode collapse but are more challenging to train while still not scaling easily  
 39 beyond tens of nodes. Recently, one-shot generation has been scaled up to graphs of hundreds of  
 40 nodes thanks to spectral conditioning [21], suggesting that good conditioning can largely benefit  
 41 graph generation. Still, the suggested training procedure is cumbersome as it involves 3 different  
 42 intertwined Generative Adversarial Networks (GANs). Finally, Variational Auto Encoders (VAE)  
 43 have also been studied to generate graphs but remain difficult to train, as the loss function needs to be  
 44 permutation invariant [22] which can necessitate an expensive graph matching step [17].

45 In contrast, the score-based models [8, 9] have the potential to provide both, a simple, stable training  
 46 objective similar to the auto-regressive models and good graph distribution coverage provided by  
 47 the one-shot models. Niu et al. [8] provided the first score-based model for graph generation by  
 48 directly using the score-based model formulation of Song and Ermon [1] and additionally accounting  
 49 for the permutation equivariance of graphs. Jo et al. [9] extended this to featured graph generation,  
 50 by formulating the problem as a system of two stochastic differential equations, one for feature  
 51 generation and one for adjacency generation. The graph and the features are then generated in  
 52 parallel. This approach provided promising results for small molecule generation. Importantly, both  
 53 contributions rely on a continuous Gaussian noise process and use a thousand denoising steps to  
 54 achieve good results, which makes for a slow graph generation.

55 As shown by Song et al. [6], score matching is tightly related to denoising diffusion probabilistic  
 56 models [3] which provide a more flexible formulation, more easily amendable for the graph generation.  
 57 In particular, for the noisy samples to remain discrete graphs, the perturbations need to be discrete.  
 58 Such discrete diffusion has been successfully used for quantized image generation [23, 24] and text  
 59 generation [25]. Diffusion using the multinomial distribution was proposed in Hoogeboom et al.  
 60 [10]. Then, Austin et al. [7] extended the previous work by Hoogeboom et al. [10], Song et al. [26]  
 61 and provided a general recipe for denoising diffusion models in discrete state-spaces which mainly  
 62 requires the specification of a doubly-stochastic Markov transition matrix  $\mathbf{Q}$  which ensures the  
 63 Markov process conserves probability mass and converges to a stationary distribution. In the next  
 64 section, we describe a formulation of this perturbation matrix  $\mathbf{Q}$  leading to the ER random graphs.

### 65 3 Discrete Diffusion for Simple Graphs

66 Diffusion models [2] are generative models based on a forward and a reverse Markov process.  
 67 The forward process, denoted  $q(\mathbf{A}_{1:T} | \mathbf{A}_0) = \prod_{t=1}^T q(\mathbf{A}_t | \mathbf{A}_{t-1})$  generates a sequence of  
 68 increasingly noisier latent variables  $\mathbf{A}_t$  from the initial sample  $\mathbf{A}_0$ , to white noise  $\mathbf{A}_T$ . Here  
 69 the sample  $\mathbf{A}_0$  and the latent variables  $\mathbf{A}_t$  are adjacency matrices. The learned reverse process  
 70  $p_\theta(\mathbf{A}_{1:T}) = p(\mathbf{A}_T) \prod_{t=1}^T q(\mathbf{A}_{t-1} | \mathbf{A}_t)$  attempts to progressively denoise the latent variable  $\mathbf{A}_t$   
 71 in order to produce samples from the desired distribution. Here we will focus on simple graphs, but the  
 72 approach can be extended in a straightforward manner to account for different edge types. We use the  
 73 model from [10] and, for convenience, adopt the representation of [7] for our discrete process.

#### 74 3.1 Forward Process

75 Let the row vector  $\mathbf{a}_t^{ij} \in \{0, 1\}^2$  be the one-hot encoding of  $i, j$  element of the adjacency matrix  $\mathbf{A}_t$ .  
 76 Here  $t \in [0, T]$  denotes the timestep of the process, where  $\mathbf{A}_0$  is a sample from the data distribution  
 77 and  $\mathbf{A}_T$  is an ER random graph. The forward process is described as repeated multiplication of each  
 78 adjacency element type row vector  $\mathbf{a}_t^{ij} = \mathbf{a}_{t-1}^{ij} \mathbf{Q}_t$  with a double stochastic matrix  $\mathbf{Q}_t$ . Note that the  
 79 forward process is independent for each edge/non-edge  $i \neq j$ . The matrix  $\mathbf{Q}_t \in \mathbb{R}^{2 \times 2}$  is modeled as

$$\mathbf{Q}_t = \begin{bmatrix} 1 - \beta_t & \beta_t \\ \beta_t & 1 - \beta_t \end{bmatrix}, \quad (1)$$

80 where  $\beta_t$  is the probability of not changing the edge state. This formulation has the advantage  
 81 to allow direct sampling at any timestep of the diffusion process without computing any previous  
 82 timesteps. Indeed the matrix  $\overline{\mathbf{Q}}_t = \prod_{i < t} \mathbf{Q}_i$  can be expressed in the form of (1) with  $\beta_t$  being

83 replaced by  $\bar{\beta}_t = \frac{1}{2} - \frac{1}{2} \prod_{i < t} (1 - 2\beta_i)$ . Eventually, we want the probability  $\bar{\beta}_t \in [0, 0.5]$  to vary  
 84 from 0 (unperturbed sample) to 0.5 (pure noise). In this contribution, we limit ourselves to symmetric  
 85 graphs and therefore only need to model the upper triangular part of the adjacency matrix. The noise  
 86 is sampled i.i.d. over all of the edges.

### 87 3.2 Reverse Process

88 To sample from the data distribution, the forward process needs to be reversed. Therefore, we need to  
 89 estimate  $q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0)$ . In our case, using the Markov property of the forward process this can be  
 90 rewritten as (see Appendix A for derivation):

$$q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0) = q(\mathbf{A}_t | \mathbf{A}_{t-1}) \frac{q(\mathbf{A}_{t-1} | \mathbf{A}_0)}{q(\mathbf{A}_t | \mathbf{A}_0)}. \quad (2)$$

91 Note that (2) is entirely defined by  $\beta_t$  and  $\bar{\beta}_t$  and  $\mathbf{A}_0$  (see Appendix A, Equation 4).

### 92 3.3 Loss

93 Diffusion models are typically trained to minimize a variational upper bound on the negative log-  
 94 likelihood. This bound can be expressed as (see Appendix C or [3, Equation 5]):

$$L_{\text{vb}}(\mathbf{A}_0) := \mathbb{E}_{q(\mathbf{A}_0)} \left[ \underbrace{D_{KL}(q(\mathbf{A}_T | \mathbf{A}_0) \| p_\theta(\mathbf{A}_T))}_{L_T} + \sum_{t=1}^T \mathbb{E}_{q(\mathbf{A}_t | \mathbf{A}_0)} \left[ \underbrace{D_{KL}(q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0) \| p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t))}_{L_t} - \underbrace{\mathbb{E}_{q(\mathbf{A}_1 | \mathbf{A}_0)} \log(p_\theta(\mathbf{A}_0 | \mathbf{A}_1))}_{L_0} \right] \right]$$

95 Practically, the model is trained to directly minimize the losses  $L_t$ , i.e. the KL divergence  
 96  $D_{KL}(q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0) \| p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t))$  by using the tractable parametrization of  $q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0)$   
 97 from (2). Note that the discrete setting of the selected noise distribution prevents training the model to  
 98 approximate the gradient of the distribution as done by score-matching graph generative models [8, 9].

99 **Parametrization of the reverse process.** While it is possible to predict the logits of  $p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t)$   
 100 in order to minimize  $L_{\text{vb}}$ , we follow [3, 10, 7] and use a network  $\text{nn}_\theta(\mathbf{A}_t)$  that predict the logits of  
 101 the distribution  $p_\theta(\mathbf{A}_0 | \mathbf{A}_t)$ . This parametrization is known to stabilize the training procedure. To  
 102 minimize  $L_{\text{vb}}$ , (2) can be used to recover  $p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t)$  from  $\mathbf{A}_0$  and  $\mathbf{A}_t$ .

103 **Alternate loss.** Many implementations of DDPMs found it beneficial to use alternative losses. For  
 104 instance, [3] derived a simplified loss function that reweights the ELBO. Hybrid losses have been  
 105 used in [27] and [7]. As shown in Appendix D, using the parametrization  $p_\theta(\mathbf{A}_0 | \mathbf{A}_t)$ , one can  
 106 express the term:  $L_t$  as  $L_t = -\log(p_\theta(\mathbf{A}_0 | \mathbf{A}_t))$ . Empirically, we found that minimizing

$$L_{\text{simple}} := -\mathbb{E}_{q(\mathbf{A}_0)} \sum_{t=1}^T \left( 1 - 2 \cdot \bar{\beta}_t + \frac{1}{T} \right) \cdot \mathbb{E}_{q(\mathbf{A}_t | \mathbf{A}_0)} \log p_\theta(\mathbf{A}_0 | \mathbf{A}_t) \quad (3)$$

107 leads to stable training and better results. Note that this loss equals the cross-entropy loss between  
 108  $\mathbf{A}_0$  and  $\text{nn}_\theta(\mathbf{A}_t)$ . The re-weighting  $1 - 2 \cdot \bar{\beta}_t + \frac{1}{T}$ , which assigns linearly more importance to the  
 109 less noisy samples, has been proposed in [23, Equation 7].

### 110 3.4 Sampling

111 For each loss, we used a specific sampling algorithm. For both approaches, we start by sampling  
 112 each edge independently from a Bernoulli distribution with probability  $p = 1/2$  (ER random graph).  
 113 Then, for the  $L_{\text{vb}}$  loss we follow Ho et al. [3] and iteratively reverse the chain by sampling Bernoulli-  
 114 sampling from  $p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t)$  until we obtain at our sample of  $p_\theta(\mathbf{A}_0 | \mathbf{A}_1)$ . For the loss function  
 115  $L_{\text{simple}}$ , we sample  $\mathbf{A}_0$  directly from  $p_\theta(\mathbf{A}_0 | \mathbf{A}_t)$  for each step  $t$  and obtain  $\mathbf{A}_{t-1}$  by sampling again  
 116 from  $q(\mathbf{A}_{t-1} | \mathbf{A}_0)$ . The two approaches are described algorithmically in Appendix E.

117 The values of  $\bar{\beta}_t$  are selected following a simple linear schedule for our reverse process [2]. We  
 118 found it works similarly well as other options such as cosine schedule [27]. Note that in this case  $\beta_t$   
 119 can be obtained from  $\bar{\beta}_t$  in a straightforward manner (see Appendix B).

## 120 4 Experiments

121 We compare our graph discrete diffusion  
 122 approach to the original score-based ap-  
 123 proach proposed by Niu et al. [8]. Models  
 124 using this original formulation are denoted  
 125 by *score*. We follow the training and evalu-  
 126 ation setup used by previous contributions  
 127 [15, 19, 8, 9]. More details can be found  
 128 in Appendix G. For evaluation, we com-  
 129 pute MMD metrics from [15] between the  
 130 generated graphs and the test set, namely,  
 131 the degree distribution, the clustering co-  
 132 efficient, and the 4-node orbit counts. To  
 133 demonstrate the efficiency of the discrete  
 134 parameterization, the discrete models only  
 135 use 32 denoising steps, while the score-  
 136 based models use 1000 denoising steps, as  
 137 originally proposed. We compare two ar-  
 138 chitectures: 1. EDP-GNN as introduced by  
 139 Niu et al. [8], and 2. a simpler and more ex-  
 140 pressive provably powerful graph network  
 141 (PPGN) [11]. See Appendix F for a more  
 142 detailed description of the architectures.

143 Table 1 shows the results for two datasets,  
 144 Community-small ( $12 \leq n \leq 20$ ) and Ego-small ( $4 \leq n \leq 18$ ), used by Niu et al. [8]. To better  
 145 compare our approach to traditional score-based graph generation, in Table 2, we additionally perform  
 146 experiments on slightly more challenging datasets with larger graphs. Namely, a stochastic-block-  
 147 model (SBM) dataset with three communities, which in total consists of ( $24 \leq n \leq 27$ ) nodes and a  
 148 planar dataset with ( $n = 60$ ) nodes. Detailed information on the datasets can be found in Appendix H.  
 149 Additional details concerning the evaluation setup are provided in Appendix G.4.

150 **Results.** In Table 1, we observe that the proposed discrete diffusion process using the  $L_{vb}$  loss and  
 151 PPGN model leads to slightly improved average MMDs over the competitors. The  $L_{simple}$  loss further  
 152 improve the result over  $L_{vb}$ . The fact that the EDP- $L_{simple}$  model has significantly lower MMD values  
 153 than the EDP-score model is a strong indication that the proposed loss and the discrete formulation  
 154 are the cause of the improvement rather than the PPGN architecture. This improvement comes with  
 155 the additional benefit that sampling is greatly accelerated (30 times) as the number of timesteps  
 156 is reduced from 1000 to 32. Table 2 shows that the proposed discrete formulation is even more  
 157 beneficial when graph size and complexity increase. The PPGN-Score even becomes infeasible to  
 158 run in this setting, due to the prohibitively expensive sampling procedure. A qualitative evaluation of  
 159 the generated graphs is performed in Appendix I. Visually, the  $L_{simple}$  loss leads to the best samples.

## 160 5 Conclusion

161 In this work, we demonstrated that discrete diffusion can increase sample quality and greatly improve  
 162 the efficiency of denoising diffusion for graph generation. While the approach was presented for  
 163 simple graphs with non-attributed edges, it could also be extended to graphs with edge attributes.

<sup>†</sup>The discrepancy with the SDE-Score<sup>†</sup> results comes from the fact that using the code provided by the  
 authors, we were unable to reproduce their results. Strangely, their code leads to good results when used with  
 our discrete formulation and  $L_{simple}$  loss improving over the result reported in their contribution.

Model	Community				Ego				Total
	Deg.	Clus.	Orb.	Avg.	Deg.	Clus.	Orb.	Avg.	
GraphRNN <sup>†</sup>	0.030	0.030	0.010	<b>0.017</b>	0.040	0.050	0.060	0.050	0.033
GNF <sup>†</sup>	0.120	0.150	0.020	0.097	0.010	0.030	0.001	0.014	0.055
EDP-Score <sup>†</sup>	0.006	0.127	0.018	0.050	0.010	0.025	0.003	<b>0.013</b>	0.031
SDE-Score <sup>†</sup>	0.045	0.086	0.007	0.046	0.021	0.024	0.007	0.017	0.032
EDP-Score <sup>†</sup>	0.016	0.810	0.110	0.320	0.04	0.064	0.005	0.037	0.178
PPGN-Score	0.081	0.237	0.284	0.200	0.019	0.049	0.005	0.025	0.113
PPGN $L_{vb}$	0.023	0.061	0.015	0.033	0.025	0.039	0.019	0.027	0.03
PPGN $L_{simple}$	0.019	0.044	0.005	0.023	0.018	0.026	0.003	0.016	<b>0.019</b>
EDP $L_{simple}$	0.024	0.04	0.012	0.026	0.019	0.031	0.017	0.022	0.024

Table 1: MMD results for the Community and the Ego datasets. All values are averaged over 5 runs with 1024 generated samples without any sub-selection. The "Total" column denotes the average MMD over all of the 6 measurements. The best results of the "Avg." and "Total" columns are shown in bold. † marks the results taken from the original papers.

Model	SBM-27				Planar-60				Total
	Deg.	Clus.	Orb.	Avg.	Deg.	Clus.	Orb.	Avg.	
EDP-Score	0.014	0.800	0.190	0.334	1.360	1.904	0.534	1.266	0.8
PPGN $L_{simple}$	0.007	0.035	0.072	<b>0.038</b>	0.029	0.039	0.036	<b>0.035</b>	<b>0.036</b>
EDP $L_{simple}$	0.046	0.184	0.064	0.098	0.017	1.928	0.785	0.910	0.504

Table 2: MMD results for the SBM-27 and the Planar-60 datasets.

## References

- [1] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019.
- [2] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.
- [3] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [4] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- [5] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, 2014*.
- [6] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations, 2021*.
- [7] Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021.
- [8] Chenhao Niu, Yang Song, Jiaming Song, Shengjia Zhao, Aditya Grover, and Stefano Ermon. Permutation invariant graph generation via score-based generative modeling. In *International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 4474–4484, Online, 26–28 Aug 2020. PMLR.
- [9] Jaehyeong Jo, Seul Lee, and Sung Ju Hwang. Score-based generative modeling of graphs via the system of stochastic differential equations. In *Proceedings of the International Conference on Machine Learning (ICML), 2022*.
- [10] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Learning categorical distributions. *Advances in Neural Information Processing Systems*, 34:12454–12465, 2021.
- [11] Haggai Maron, Heli Ben-Hamu, Hadar Serviansky, and Yaron Lipman. Provably powerful graph networks. In *Advances in Neural Information Processing Systems*, pages 2156–2167, 2019.
- [12] Paul Erdos, Alfréd Rényi, et al. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci.*, 5(1):17–60, 1960.
- [13] Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First steps. *Social networks*, 5(2):109–137, 1983.
- [14] Justin Eldridge, Mikhail Belkin, and Yusu Wang. Graphons, mergeons, and so on! In *Advances in Neural Information Processing Systems*, pages 2307–2315, 2016.
- [15] Jiaxuan You, Rex Ying, Xiang Ren, William L Hamilton, and Jure Leskovec. Graphrnn: Generating realistic graphs with deep auto-regressive models. In *ICML*, 2018.
- [16] Renjie Liao, Yujia Li, Yang Song, Shenlong Wang, Will Hamilton, David K Duvenaud, Raquel Urtasun, and Richard Zemel. Efficient graph generation with graph recurrent attention networks. In *Advances in Neural Information Processing Systems*, pages 4255–4265, 2019.
- [17] Martin Simonovsky and Nikos Komodakis. Graphvae: Towards generation of small graphs using variational autoencoders. In *International conference on artificial neural networks*, pages 412–422. Springer, 2018.

- 210 [18] Nicola De Cao and Thomas Kipf. MolGAN: An implicit generative model for small molecular  
211 graphs. *ICML 2018 workshop on Theoretical Foundations and Applications of Deep Generative*  
212 *Models*, 2018.
- 213 [19] Jenny Liu, Aviral Kumar, Jimmy Ba, Jamie Kiros, and Kevin Swersky. Graph normalizing  
214 flows. *Advances in Neural Information Processing Systems*, 32:13578–13588, 2019.
- 215 [20] Igor Krawczuk, Pedro Abranches, Andreas Loukas, and Volkan Cevher. Gg-gan: A geometric  
216 graph generative adversarial network. 2020.
- 217 [21] Karolis Martinkus, Andreas Loukas, Nathanaël Perraudin, and Roger Wattenhofer. Spectre:  
218 Spectral conditioning helps to overcome the expressivity limits of one-shot graph generators. In  
219 *Proceedings of the International Conference on Machine Learning (ICML)*, 2022.
- 220 [22] Clement Vignac and Pascal Frossard. Top-n: Equivariant set and graph generation without  
221 exchangeability. In *International Conference on Learning Representations*, 2022.
- 222 [23] Sam Bond-Taylor, Peter Hesse, Hiroshi Sasaki, Toby P. Breckon, and Chris G. Willcocks.  
223 Unleashing transformers: Parallel token prediction with discrete absorbing diffusion for fast  
224 high-resolution image generation from vector-quantized codes. In *European Conference on*  
225 *Computer Vision (ECCV)*, 2022.
- 226 [24] Patrick Esser, Robin Rombach, Andreas Blattmann, and Bjorn Ommer. Imagebart: Bidirectional  
227 context with multinomial diffusion for autoregressive image synthesis. *Advances in Neural*  
228 *Information Processing Systems*, 34:3518–3532, 2021.
- 229 [25] Emiel Hoogeboom, Alexey A. Gritsenko, Jasmijn Bastings, Ben Poole, Rianne van den Berg,  
230 and Tim Salimans. Autoregressive diffusion models. In *International Conference on Learning*  
231 *Representations*, 2022.
- 232 [26] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In  
233 *International Conference on Learning Representations*, 2020.
- 234 [27] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic  
235 models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.
- 236 [28] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural  
237 networks? In *International Conference on Learning Representations*, 2019.
- 238 [29] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-  
239 Rad. Collective classification in network data articles. *AI Magazine*, 29:93–106, 09 2008. doi:  
240 10.1609/aimag.v29i3.2157.
- 241 [30] Der-Tsai Lee and Bruce J Schachter. Two algorithms for constructing a delaunay triangulation.  
242 *International Journal of Computer & Information Sciences*, 9(3):219–242, 1980.

## 243 A Reverse Process Derivations

244 In this appendix, we provide the derivation of the reverse probability  $q(\mathbf{A}_{t-1}|\mathbf{A}_t, \mathbf{A}_0)$ . Using the  
 245 Bayes rule, we obtain

$$\begin{aligned} q(\mathbf{A}_{t-1}|\mathbf{A}_t, \mathbf{A}_0) &= \frac{q(\mathbf{A}_t | \mathbf{A}_{t-1}, \mathbf{A}_0) \cdot q(\mathbf{A}_{t-1}, \mathbf{A}_0)}{q(\mathbf{A}_t, \mathbf{A}_0)} \\ &= \frac{q(\mathbf{A}_t | \mathbf{A}_{t-1}) \cdot q(\mathbf{A}_{t-1} | \mathbf{A}_0)q(\mathbf{A}_0)}{q(\mathbf{A}_t | \mathbf{A}_0) \cdot q(\mathbf{A}_0)} \\ &= q(\mathbf{A}_t | \mathbf{A}_{t-1}) \cdot \frac{q(\mathbf{A}_{t-1} | \mathbf{A}_0)}{q(\mathbf{A}_t | \mathbf{A}_0)}, \end{aligned}$$

246 where we use the fact that  $q(\mathbf{A}_t | \mathbf{A}_{t-1}, \mathbf{A}_0) = q(\mathbf{A}_t | \mathbf{A}_{t-1})$  since  $\mathbf{A}_t$  is independent of  $\mathbf{A}_0$  given  
 247  $\mathbf{A}_{t-1}$ .

248 This reverse probability is entirely defined with  $\beta_t$  and  $\bar{\beta}_t$ . For the  $i, j$  element of  $\mathbf{A}$  (denoted  $\mathbf{A}^{ij}$ ),  
 249 we obtain:

$$q(\mathbf{A}_{t-1}^{ij} = 1 | \mathbf{A}_t^{ij}, \mathbf{A}_0^{ij}) = \begin{cases} (1 - \beta_t) \cdot \frac{(1 - \bar{\beta}_{t-1})}{1 - \bar{\beta}_t}, & \text{if } \mathbf{A}_t^{ij} = 1, \mathbf{A}_0^{ij} = 1 \\ (1 - \beta_t) \cdot \frac{\bar{\beta}_{t-1}}{\bar{\beta}_t}, & \text{if } \mathbf{A}_t^{ij} = 1, \mathbf{A}_0^{ij} = 0 \\ \beta_t \cdot \frac{(1 - \bar{\beta}_{t-1})}{\bar{\beta}_t}, & \text{if } \mathbf{A}_t^{ij} = 0, \mathbf{A}_0^{ij} = 1 \\ \beta_t \cdot \frac{\bar{\beta}_{t-1}}{1 - \bar{\beta}_t}, & \text{if } \mathbf{A}_t^{ij} = 0, \mathbf{A}_0^{ij} = 0 \end{cases} \quad (4)$$

## 250 B Conversion of $\bar{\beta}_t$ to $\beta_t$

251 The selected linear schedule provides us with the values of  $\bar{\beta}_t$ . In this appendix, we compute  
 252 an expression for  $\beta_t$  from  $\bar{\beta}_t$ , which allows us easy computation of (2). By definition, we have  
 253  $\bar{\mathbf{Q}}_t = \bar{\mathbf{Q}}_{t-1} \mathbf{Q}_t$  which is equivalent to

$$\begin{pmatrix} 1 - \bar{\beta}_{t-1} & \bar{\beta}_{t-1} \\ \bar{\beta}_{t-1} & 1 - \bar{\beta}_{t-1} \end{pmatrix} \begin{pmatrix} 1 - \beta_t & \beta_t \\ \beta_t & 1 - \beta_t \end{pmatrix} = \begin{pmatrix} 1 - \bar{\beta}_t & \bar{\beta}_t \\ \bar{\beta}_t & 1 - \bar{\beta}_t \end{pmatrix}$$

254 Let us select the first row and first column equality. We obtain the following equation

$$(1 - \bar{\beta}_{t-1})(1 - \beta_t) + \bar{\beta}_{t-1}\beta_t = 1 - \bar{\beta}_t,$$

255 which, after some arithmetic, provides us with the desired answer

$$\beta_t = \frac{\bar{\beta}_{t-1} - \bar{\beta}_t}{2\bar{\beta}_{t-1} - 1}.$$

## 256 C ELBO derivation

257 The general Evidence Lower Bound (ELBO) formula states that

$$\log(p_\theta(x)) \geq \mathbb{E}_{z \sim q} \left[ \log \left( \frac{p(x, z)}{q(z)} \right) \right]$$

258 for any distribution  $q$  and latent  $z$ . In our case, we use  $\mathbf{A}_{1:T}$  as a latent variable and obtain

$$-\log(p_\theta(\mathbf{A}_0)) \leq \mathbb{E}_{\mathbf{A}_{1:T} \sim q(\mathbf{A}_{1:T}|\mathbf{A}_0)} \left[ \log \left( \frac{p_\theta(\mathbf{A}_{0:T})}{q(\mathbf{A}_{1:T} | \mathbf{A}_0)} \right) \right] := L_{\text{vb}}(\mathbf{A}_0)$$

259 We use  $L_{\text{vb}} = \mathbb{E}[L_{\text{vb}}(\mathbf{A}_0)]$  and obtain

$$\begin{aligned}
L_{\text{vb}} &= \mathbb{E}_{q(\mathbf{A}_{0:T})} \left[ -\log \left( \frac{p_\theta(\mathbf{A}_{0:T})}{q(\mathbf{A}_{1:T} | \mathbf{A}_0)} \right) \right] \\
&= \mathbb{E}_q \left[ -\log(p_\theta(\mathbf{A}_T)) - \sum_{t=1}^T \log \left( \frac{p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t)}{q(\mathbf{A}_t | \mathbf{A}_{t-1})} \right) \right] \\
&= \mathbb{E}_q \left[ -\log(p_\theta(\mathbf{A}_T)) - \sum_{t=2}^T \log \left( \frac{p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t)}{q(\mathbf{A}_t | \mathbf{A}_{t-1})} \right) - \log \left( \frac{p_\theta(\mathbf{A}_0 | \mathbf{A}_1)}{q(\mathbf{A}_1 | \mathbf{A}_0)} \right) \right] \\
&= \mathbb{E}_q \left[ -\log(p_\theta(\mathbf{A}_T)) - \sum_{t=2}^T \log \left( \frac{p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t)}{q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0)} \cdot \frac{q(\mathbf{A}_{t-1} | \mathbf{A}_0)}{q(\mathbf{A}_t | \mathbf{A}_0)} \right) - \log \left( \frac{p_\theta(\mathbf{A}_0 | \mathbf{A}_1)}{q(\mathbf{A}_1 | \mathbf{A}_0)} \right) \right] \\
&= \mathbb{E}_q \left[ -\log \left( \frac{p_\theta(\mathbf{A}_T)}{q(\mathbf{A}_T | \mathbf{A}_0)} \right) - \sum_{t=2}^T \log \left( \frac{p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t)}{q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0)} \right) - \log(p_\theta(\mathbf{A}_0 | \mathbf{A}_1)) \right] \\
&= \mathbb{E}_{\mathbb{E}_{q(\mathbf{A}_0)}} \left[ D_{KL}(q(\mathbf{A}_T | \mathbf{A}_0) \| p_\theta(\mathbf{A}_T)) + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{A}_t | \mathbf{A}_0)} D_{KL}(q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0) \| p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t)) \right. \\
&\quad \left. - \mathbb{E}_{q(\mathbf{A}_1 | \mathbf{A}_0)} \log(p_\theta(\mathbf{A}_0 | \mathbf{A}_1)) \right] \tag{5}
\end{aligned}$$

260 where (5) follows from

$$\begin{aligned}
q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0) &= \frac{q(\mathbf{A}_t | \mathbf{A}_{t-1}, \mathbf{A}_0) q(\mathbf{A}_{t-1}, \mathbf{A}_0)}{q(\mathbf{A}_t, \mathbf{A}_0)} \\
&= \frac{q(\mathbf{A}_t | \mathbf{A}_{t-1}) q(\mathbf{A}_{t-1} | \mathbf{A}_0)}{q(\mathbf{A}_t | \mathbf{A}_0)}.
\end{aligned}$$

## 261 D Simple Loss

262 Using the parametrization  $p_\theta(\mathbf{A}_0 | \mathbf{A}_t)$ , we can simplify the KL divergenc of the term  $L_t$ .

$$\begin{aligned}
D_{KL}(q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0) \| p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t)) &= \mathbb{E}_{q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0)} \left[ -\log \left( \frac{p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t)}{q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0)} \right) \right] \\
&= \mathbb{E}_{q(\mathbf{A}_{t-1} | \mathbf{A}_t, \mathbf{A}_0)} [-\log(p_\theta(\mathbf{A}_0 | \mathbf{A}_t))] \\
&= -\log(p_\theta(\mathbf{A}_0 | \mathbf{A}_t))
\end{aligned}$$

263 We note that this term corresponds to the cross-entropy of the distribution  $p_\theta(\mathbf{A}_0 | \mathbf{A}_t)$  with the  
264 ground truth of  $\mathbf{A}_0$ .

## 265 E Sampling Algorithms

266 Here in Algorithms 1 and 2 we provide an algorithmic description of the two sampling approaches  
267 described in Section 3.4. Here  $\mathcal{B}_{p=1/2}$  denotes the Bernoulli distribution with parameter  $p = 1/2$ ,  
268 which corresponds to the Erdős–Rényi random graph model.

---

### Algorithm 1 Sampling for $L_{\text{vb}}$

---

1:  $\forall i, j | i > j: \mathbf{A}_T^{ij} \sim \mathcal{B}_{p=1/2}$   
2: **for**  $t = T, \dots, 1$  **do**  
3:     Compute  $p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t)$   
4:      $\mathbf{A}_{t-1} \sim p_\theta(\mathbf{A}_{t-1} | \mathbf{A}_t)$   
5: **end for**

---



---

### Algorithm 2 Sampling for $L_{\text{simple}}$

---

1:  $\forall i, j | i > j: \mathbf{A}_T^{ij} \sim \mathcal{B}_{p=1/2}$   
2: **for**  $t = T, \dots, 1$  **do**  
3:      $\tilde{\mathbf{A}}_0 \sim p_\theta(\mathbf{A}_0 | \mathbf{A}_t)$   
4:      $\mathbf{A}_{t-1} \sim q(\mathbf{A}_{t-1} | \tilde{\mathbf{A}}_0)$   
5: **end for**

---



270 **F Models**

271 **F.1 Edgewise Dense Prediction Graph Neural Network (EDP-GNN)**

272 The EDP-GNN model introduced by Niu et al. [8] extends GIN [28] to work with multi-channel  
 273 adjacency matrices. This means that a GIN graph neural network is run on multiple different adjacency  
 274 matrices (channels) and the different outputs are concatenated to produce new node embeddings:

$$\mathbf{X}_c^{(k+1)'} = \tilde{\mathbf{A}}_c^{(k)} \mathbf{X}^{(k)} + (1 + \epsilon) \mathbf{X}^{(k)},$$

275

$$\mathbf{X}^{(k+1)} = \text{Concat}(\mathbf{X}_c^{(k+1)'} \text{ for } c \in \{1, \dots, C^{(k+1)}\}),$$

276 where  $\mathbf{X} \in \mathbb{R}^{n \times h}$  is the node embedding matrix with hidden dimension  $h$  and  $C^{(k)}$  is the number of  
 277 channels in the input multi-channel adjacency matrix  $\tilde{\mathbf{A}}^{(k)} \in \mathbb{R}^{C^{(k)} \times n \times n}$ , at layer  $k$ . The adjacency  
 278 matrices for the next layer are produced using the node embeddings:

$$\tilde{\mathbf{A}}_{.,i,j}^{(k+1)} = \text{MLP}(\tilde{\mathbf{A}}_{.,i,j}^{(k)}, \mathbf{X}_i, \mathbf{X}_j).$$

279 For the first layer, EDP-GNN computes two adjacency matrix  $\tilde{\mathbf{A}}^{(0)}$  channels, original input adjacency  
 280  $\mathbf{A}$  and its inversion  $\mathbf{1}\mathbf{1}^T - \mathbf{A}$ . For node features, node degrees are used  $\mathbf{X}^{(0)} = \sum_i \mathbf{A}_i$ .

281 To produce the final outputs, outputs of all intermediary layers are concatenated:

$$\tilde{\mathbf{A}} = \text{MLP}_{\text{out}}(\text{Concat}(\tilde{\mathbf{A}}^{(k)} \text{ for } k \in \{1, \dots, K\})).$$

282 The final layer always has only one output channel, such that  $\mathbf{A}_{(t)} = \text{EDP-GNN}(\mathbf{A}_{(t-1)})$ .

283 To condition the model on the given noise level  $\bar{\beta}_t$ , noise-level-dependent scale and bias parameters  
 284  $\alpha_t$  and  $\gamma_t$  are introduced to each layer  $f$  of every MLP:

$$f(\tilde{\mathbf{A}}_{.,i,j}) = \text{activation}((\mathbf{W} \tilde{\mathbf{A}}_{.,i,j} + \mathbf{b})\alpha_t + \gamma_t).$$

285 **F.2 Provably Powerful Graph Network (PPGN)**

286 The input to the PPGN model used is the adjacency matrix  $\mathbf{A}_t$  concatenated with the diagonal matrix  
 287  $\bar{\beta}_t \cdot \mathbf{I}$ , resulting in an input tensor  $\mathbf{X}_{in} \in \mathbb{R}^{n \times n \times 2}$ . The output tensor is  $\mathbf{X}_{out} \in \mathbb{R}^{n \times n \times 1}$ , where  
 288 each  $[\mathbf{X}_{out}]_{ij}$  represents  $p([\mathbf{A}_0]_{ij} \mid [\mathbf{A}_t]_{ij})$ .

289 Our PPGN implementation, which closely follows Maron et al. [11] is structured as follows:

290 Let  $\mathbf{P}$  denote the PPGN model, then

$$\mathbf{P}(\mathbf{X}_{in}) := (l_{\text{out}} \circ C)(\mathbf{X}_{in}) \tag{6}$$

291

$$C : \mathbb{R}^{n \times n \times 2} \rightarrow \mathbb{R}^{n \times n \times (d \cdot h)} \tag{7}$$

292

$$C(\mathbf{X}_{in}) := \text{Concat}((B_d \circ \dots \circ B_1)(\mathbf{X}_{in}), (B_{d-1} \circ \dots \circ B_1)(\mathbf{X}_{in}), \dots, B_1(\mathbf{X}_{in})) \tag{8}$$

293 The set  $\{B_1, \dots, B_d\}$  is a set of  $d$  different powerful layers implemented as proposed by Maron et al.  
 294 [11]. We let the input run through different amounts of these powerful layers and concatenate their  
 295 respective outputs to one tensor of size  $n \times n \times (d \cdot h)$ . These powerful layers are functions of size:

$$\forall B_i \in \{B_2, \dots, B_d\}, B_i : \mathbb{R}^{n \times n \times h} \rightarrow \mathbb{R}^{n \times n \times h} \tag{9}$$

296

$$B_1 : \mathbb{R}^{n \times n \times 1} \rightarrow \mathbb{R}^{n \times n \times h}. \tag{10}$$

297 Finally, we use an MLP 2 to reduce the dimensionality of each matrix element down to 1, so that we  
 298 can treat the output as an adjacency matrix.

$$l_{\text{out}} : \mathbb{R}^{d \cdot h} \rightarrow \mathbb{R}^1, \tag{11}$$

299 where  $l_{\text{out}}$  is applied to each element  $[C(\mathbf{X}_{in})]_{i,j, \cdot}$  of the tensor  $C(\mathbf{X}_{in})$  over all its  $d \cdot h$  channels. It  
 300 is used to reduce the number of channels down to a single one which represents  $p(\mathbf{A}_0 \mid \mathbf{A}_t)$ .

## 301 G Training Setup

### 302 G.1 EDP-GNN

303 The model training setup and hyperparameters used for the EDP-GNN were directly taken from [8].  
304 We used 4 message-passing steps for each GIN, then stacked 5 EDP-GNN layers, for which the  
305 maximum number of channels is always set to 4 and the maximum number of node features is 16.  
306 We use 32 denoising steps for all datasets besides Planar-60, where we used 256. Opposed to 6 noise  
307 levels with 1000 sample steps per level as in the Score-based approach.

### 308 G.2 PPGN

309 The PPGN model we used for the Ego-small, Community-small, and SBM-27 datasets consist of  
310 6 layers  $\{B_1, \dots, B_6\}$ . After each powerful layer, we apply an instance normalization. The hidden  
311 dimension was set to 16. For the Planar-60 dataset, we have used 8 layers and a hidden dimension of  
312 128. We used a batch size of 64 for all datasets and used the Adam optimizer with parameters chosen  
313 as follows: learning rate is 0.001, betas are (0.9, 0.999) and weight decay is 0.999.

### 314 G.3 Model Selection

315 We performed a simple model selection where the model which achieves the best training loss is  
316 saved and used to generate graphs for testing. We also investigated the use of a validation split and  
317 computation of MMD scores versus this validation split for model selection, but we did not find this  
318 to produce better results while adding considerable computational overhead.

### 319 G.4 Additional Details on Experimental Setup

320 Here we provide some details concerning the experimental setup for the results in Tables 1 and 2.

321 **Details for MMD results in Table 1:** From the original paper Niu et al. [8], we are unsure if the  
322 GNF, GraphRNN, and EDP-Score model selection were used or not. The SDE-Score results in the  
323 first section are sampled after training for 5000 epochs and no model selection was used. Due to the  
324 compute limitations on the PPGN model, the results for PPGN  $L_{vb}$  are taken after epoch 900 instead  
325 of 5000, as results for SDE-Score and EDP-Score have been. The results for PPGN  $L_{simple}$  and EDP  
326  $L_{simple}$  were trained for 2500 epochs.

327 **Details for MMD results in Table 2:** All results using the EDP-GNN model are trained until epoch  
328 5000 and the PPGN implementation was trained until epoch 2500.

## 329 H Datasets

330 In this appendix, we describe the 4 datasets used in our experiments.

331 **Ego-small:** This dataset is composed of 200 graphs of 4-18 nodes from the Citeseer network (Sen  
332 et al. [29]). The dataset is available in the repository<sup>2</sup> of Niu et al. [8].

333 **Community-small:** This dataset consists of 100 graphs from 12 to 20 nodes. The graphs are  
334 generated in two steps. First two communities of equal size are generated using the Erdos-Rényi  
335 model [12] with parameter  $p = 0.7$ . Then edges are randomly added between the nodes of the two  
336 communities with a probability  $p = 0.05$ . The dataset is directly taken from the repository of Niu  
337 et al. [8].

338 **SBM-27:** This dataset consists of 200 graphs with 24 to 27 nodes generated using the Stochastic-  
339 Block-Model (SBM) with three communities. We use the implementation provided by Martinkus  
340 et al. [21]. The parameters used are  $p_{intra} = 0.85$ ,  $p_{inter}=0.046875$ , where  $p_{intra}$  stands for the  
341 intra-community (i.e. for node within the same community) edge probability and  $p_{inter}$  stands for the

---

<sup>2</sup><https://github.com/ermongroup/GraphScoreMatching>

342 inter-community (i.e. for nodes from different community) edge probability. The number of nodes  
343 for the 3 communities is randomly drawn from  $\{7, 8, 9\}$ . In expectation, these parameters generate 3  
344 edges between each pair of communities.

345 **Planar-60:** This dataset consists of 200 randomly generated planar graphs of 60 nodes. We use  
346 the implementation provided by Martinkus et al. [21]. To generate a graph, 60 points are first  
347 random uniformly sampled on the  $[0, 1]^2$  plane. Then the graph is generated by applying Delaunay  
348 triangulation to these points [30].

## 349 **I Visualization of Sampled Graphs**

350 In the following pages, we provide a visual comparison of graphs generated by the different models.

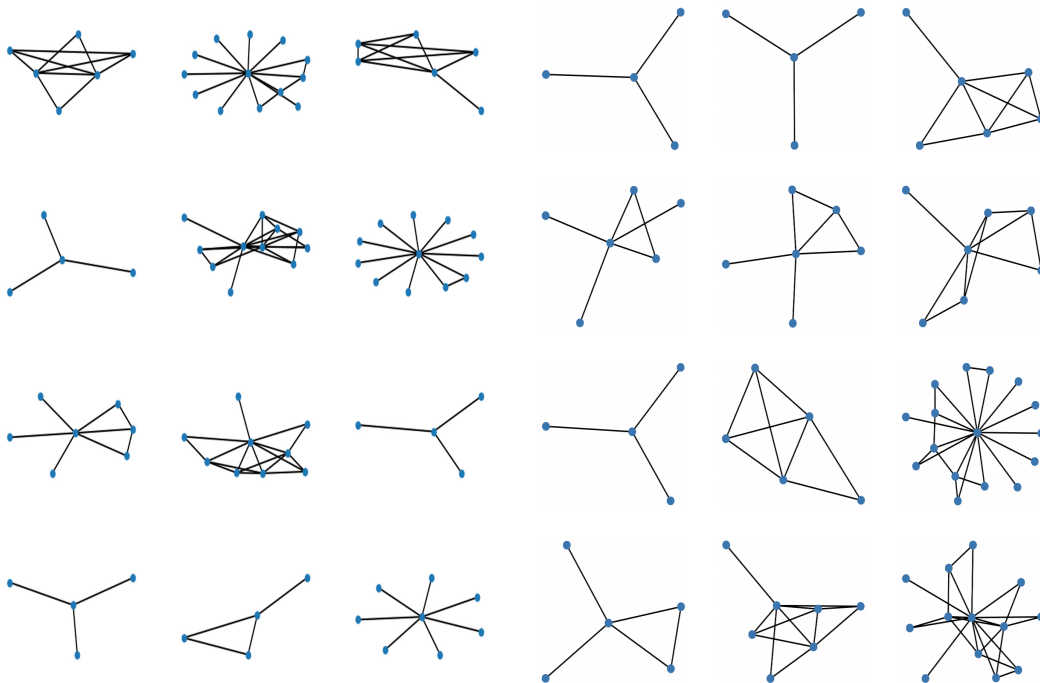


Figure 1: Sample graphs from the training set of Ego-small dataset.

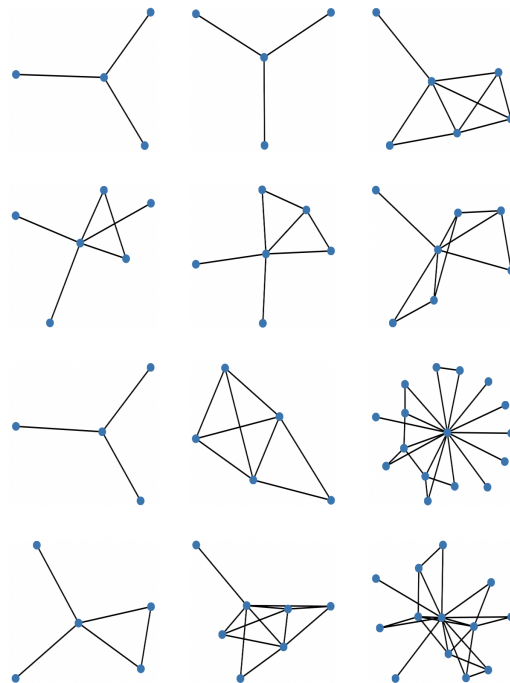


Figure 2: Sample graphs generated with the model EDP-Score [8] for the Ego-small dataset.

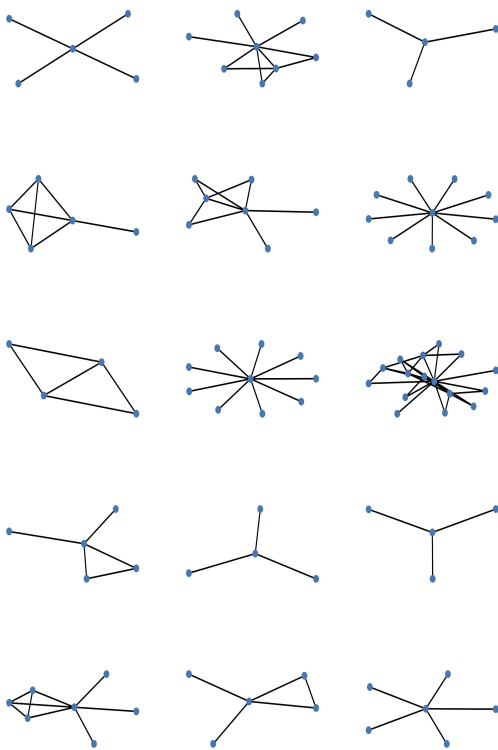


Figure 3: Sample graphs generated with the PPGN  $L_{vb}$  model for the Ego-small dataset.

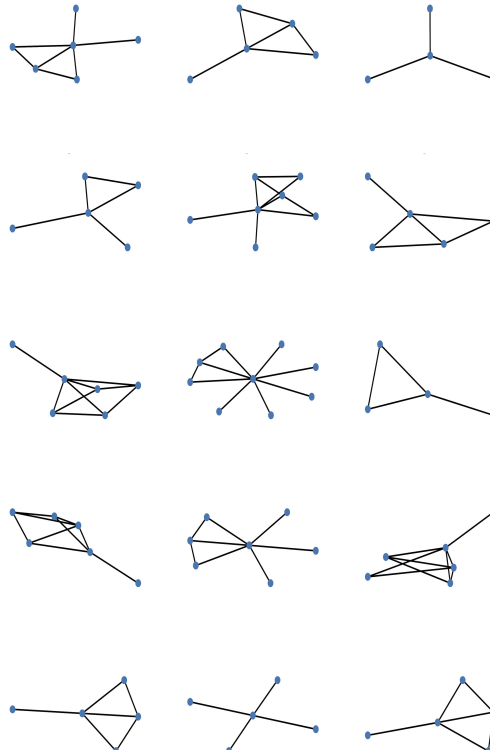


Figure 4: Sample graphs generated with the EDP  $L_{simple}$  model for the Ego-small dataset.

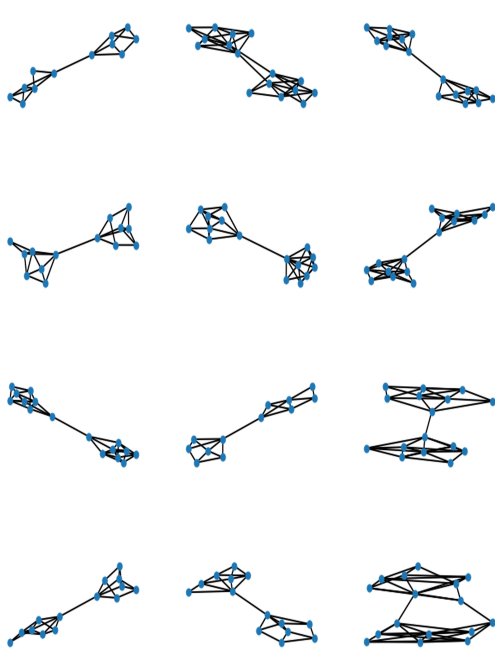


Figure 5: Sample graphs from the training set of the Community dataset

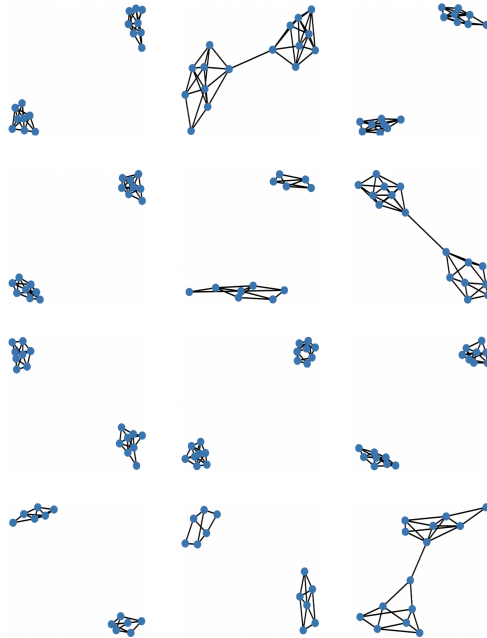


Figure 6: Sample graphs generated with the model EDP-Score [8] for the Community dataset.

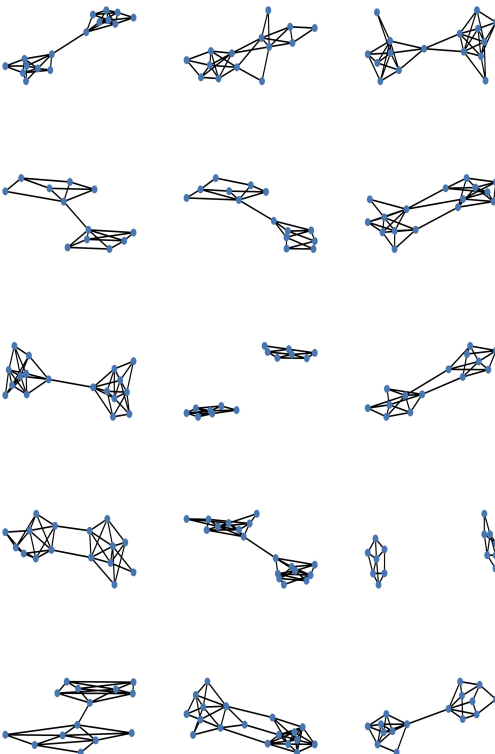


Figure 7: Sample graphs generated with the PPGN  $L_{vb}$  model for the Community dataset.

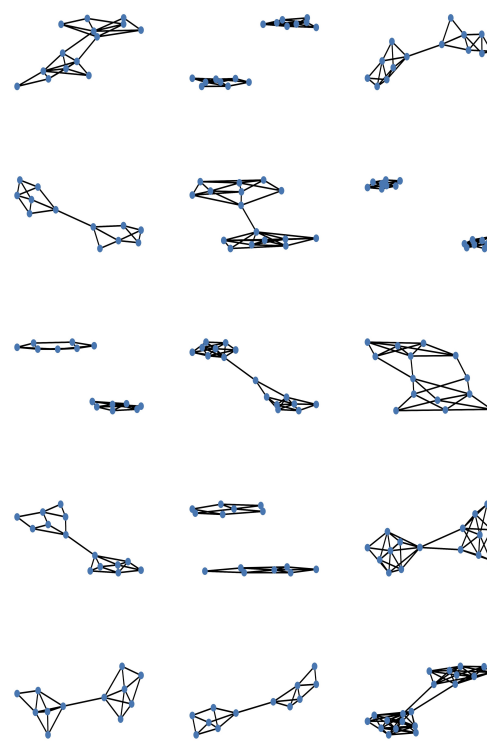


Figure 8: Sample graphs generated with the EDP  $L_{simple}$  model for the Community dataset.

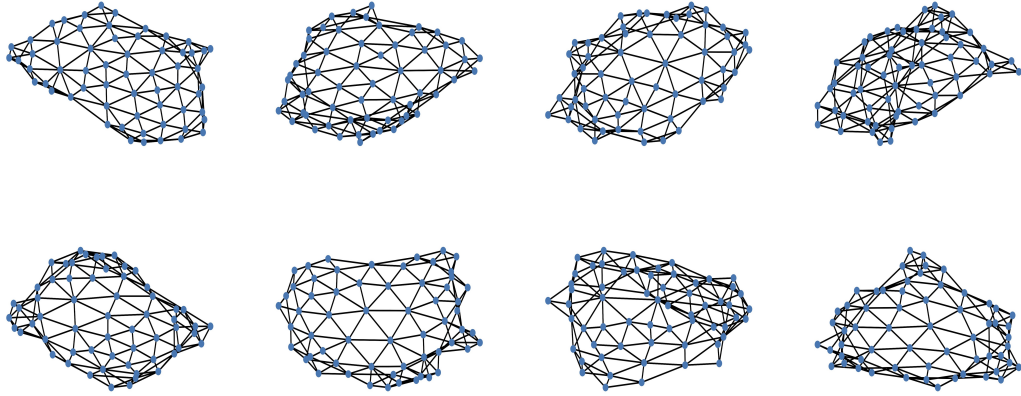


Figure 9: Sample graphs from the training set of the Planar-60 dataset.

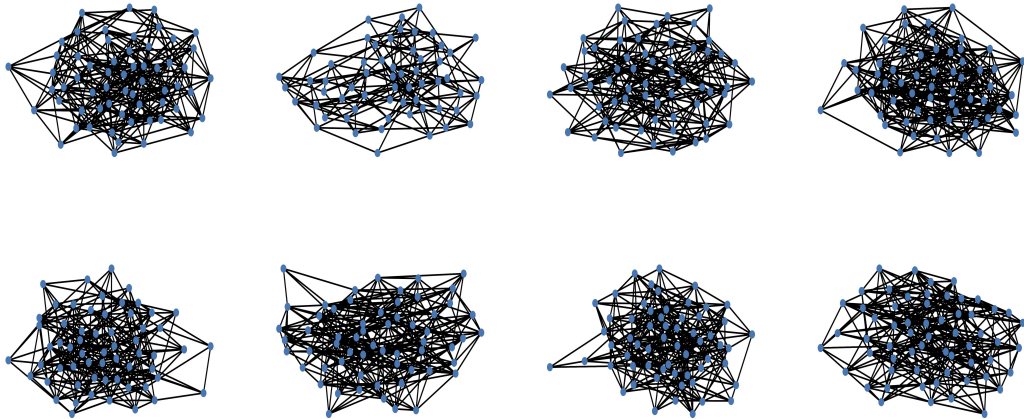


Figure 10: Sample graphs generated with the model EDP-Score [8] for the Planar-60 dataset.

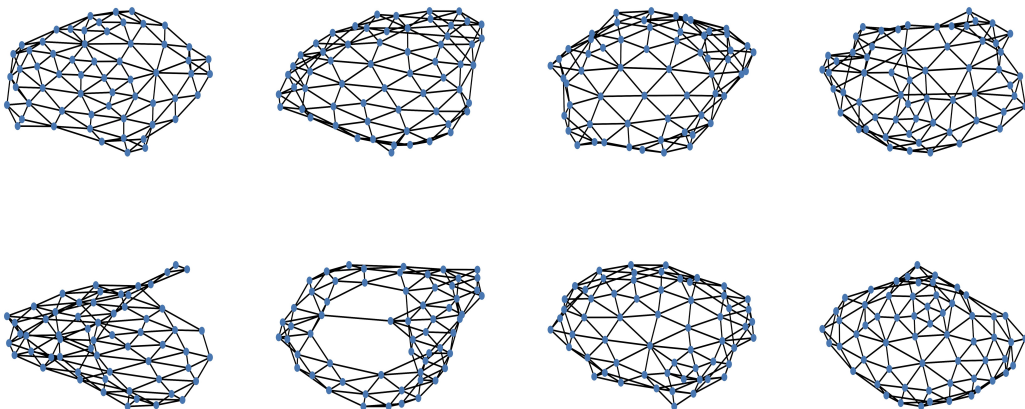


Figure 11: Sample graphs generated with the PPGN  $L_{simple}$  model for the Planar-60 dataset.

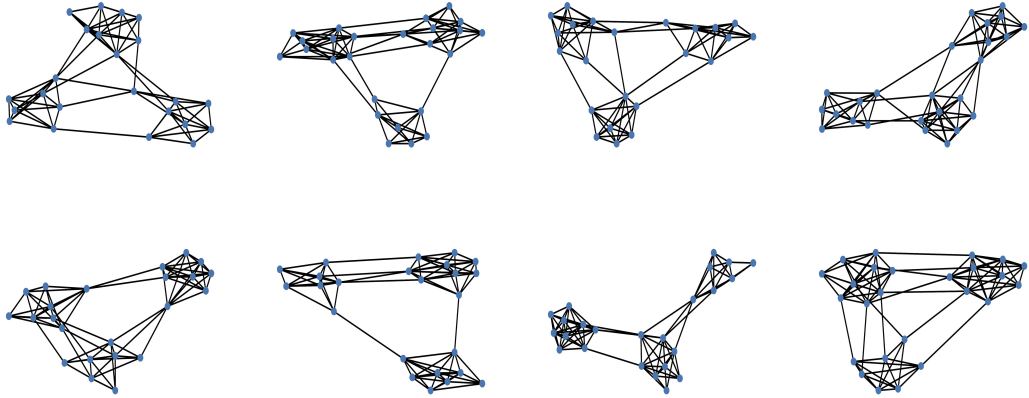


Figure 12: Sample graphs from the training set of the SBM-27 dataset.

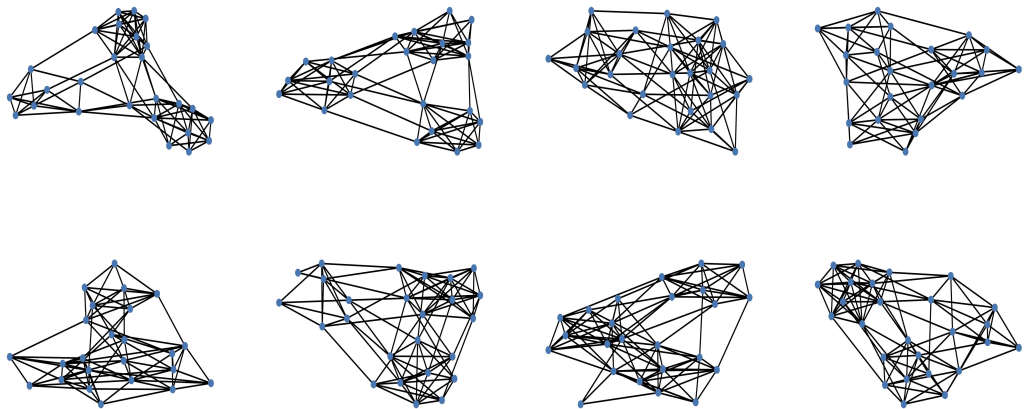


Figure 13: Sample graphs generated with the model EDP-Score [8] for the SBM-27 dataset.

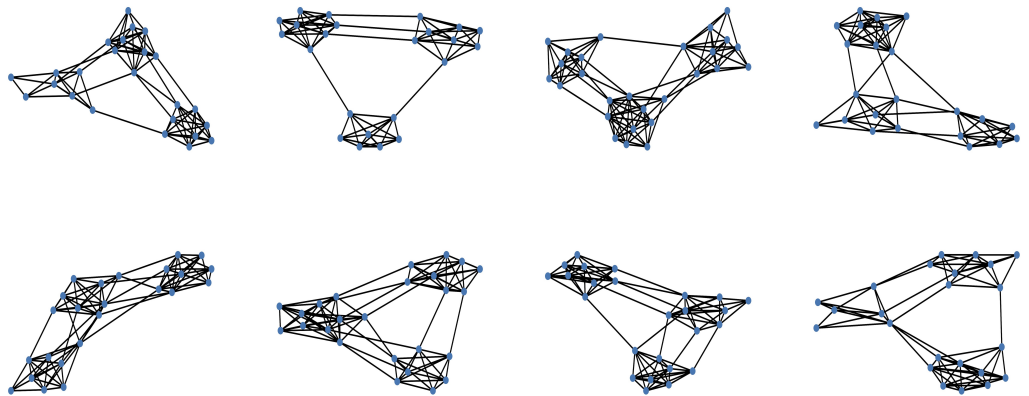


Figure 14: Sample graphs generated with the PPGN  $L_{simple}$  model for the SBM-27 dataset.