
Bures-Wasserstein Flow Matching for Graph Generation

Keyue Jiang* **Aaron Cui** **Laura Toni**
University College London, London, UK

Xiaowen Dong
University of Oxford, Oxford, UK

Abstract

Graph generation has emerged as a critical task in fields ranging from molecule design to drug discovery. Contemporary approaches, notably diffusion and flow-based models, have achieved solid graph generative performance through constructing a probability path that interpolates between a reference distribution and the data distribution. However, these methods typically model the evolution of individual nodes and edges independently and use linear interpolations to build the path assuming that the data lie in Euclidean space. We show that this is suboptimal given the intrinsic non-Euclidean structure and interconnected patterns of graphs, and it poses risks to the sampling convergence. To build a better probability path, we model the joint evolution of the nodes and edges by representing graphs as connected systems parameterized by Markov random fields (MRF). We then leverage the optimal transport displacement between MRF objects to design the probability path for graph generation. Based on this, we introduce BWFlow, a flow-matching framework for graph generation that respects the underlying geometry of graphs and provides smooth velocities in the probability path. The novel framework can be adapted to both continuous and discrete flow-matching algorithms. Experimental evaluations in plain graph generation and 2D/3D molecule generation validate the effectiveness of BWFlow in graph generation with competitive performance, stable training, and guaranteed sampling convergence.

1 Introduction

Thanks to the capability of graphs in representing complex relationships, graph generation [65, 35] has become an essential task in various fields such as protein design [26], drug discovery [6], and social network analysis [31]. Among contemporary generative models, diffusion and flow models have emerged as two compelling approaches for their ability to achieve state-of-the-art performance in graph generation [41, 55, 18, 45, 25]. In particular, these generative models can be unified under the framework of stochastic interpolation [1], which consists of four procedures [34]: 1) Drawing samples from the reference (source) distribution $p_0(\cdot)$ and/or the data (target) distribution $p_1(\cdot)$ for training set assembly; 2) Constructing a time-continuous probability path $p_t(\cdot)$, $0 \leq t \leq 1$ interpolating between p_0 and p_1 ; 3) Training a model to reconstruct the probability path by either approximating the score function or velocity fields (ratio matrix in the discrete case); and 4) sampling from p_0 and transforming it through the learned probability path to get samples that approximately follow p_1 .

A core challenge in this framework is constructing the probability path p_t . Existing text and image generative models, operating either in the continuous [24, 49, 33, 36] or discrete [8, 51, 9, 19, 40] space, typically rely on linear interpolation between source and target distributions to construct the path. Graph generation models, including diffusion [41, 55, 23, 60, 48] and flow-based models [18, 45, 25], inherit this design by modeling every single node and edge independently and linearly build paths in the disjoint space. However, this approach is inefficient because it neglects the strong

*keyue.jiang.18@ucl.ac.uk

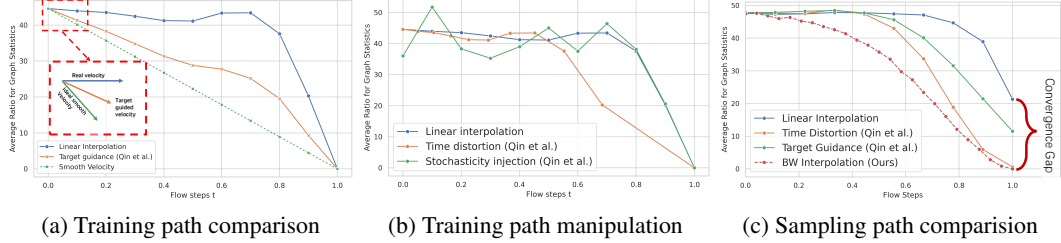


Figure 1: Probability path visualization. Since the probability is intractable, the average maximum mean discrepancy ratio (y-axis) of graph statistics between interpolants and the data points is used as a proxy for the probability. Lower means closer to the data distribution (details in Appendix I.3).

interactions and relational structure inherent in graphs, i.e., the significance of a node heavily depends on the configuration of its neighbors. While empirical success have been achieved via fine-grained searching on the training and sampling design [45] such as target guidance and time distortion, we argue that there remains a fundamental issue of the linear probability path construction, and these strategies only mitigate the problem by manipulating the probability path.

Motivation examples. The blue line in Fig. 1a illustrates the probability path evolution through linear interpolation in plain graph generation, where the probability path remains flat until $t \approx 0.8$ before sharply dropping. This pattern provides a poor velocity² estimation, as ideally the velocity field should be smooth and consistently pointing to the data distribution as the green line in Fig. 1a [29]. The resultant velocity will make the sampling difficult to converge to the data distribution, as shown in Fig. 1c, where the termination point of the blue curve remains high. Though not explicitly mentioned, the superior performance achieved by Qin et al. [45] is partially attributed to an extensive design space search to manipulate the path for smoother velocity estimation. The techniques they used, including target guidance, time distortion, and stochasticity injection, are conceptually visualized in Fig. 1a and 1b with discussions in Appendix F.1. More motivating examples in Appendix I.2.

The limitations. The above examples reveal fundamental issues of the probability path construction in graph generation, attributable to two primary reasons: 1) The assumption of independence between nodes/edges and a linear interpolation in the locally disjoint space fails to capture the global co-evolution of graph components and properties such as community structure or spectrum [22]. This potentially causes a sudden transition from reference to data distribution, which yields non-smooth probability paths. 2) The linear interpolation is derived through the optimal transport (OT) displacement [54] between distributions residing in the Euclidean space [34]. However, linear interpolation would stray away from the true data manifold when the underlying space is non-Euclidean [11, 29]. Since graphs naturally inhabit non-Euclidean geometries, linearly interpolating the nodes/edges neither guarantees an OT displacement nor respects the underlying geometry, making the constructed probability path suboptimal or even deviate from the valid graph domain.

Proposed solution. To address these limitations, we draw on statistical relational learning and model graphs using Markov Random Fields (MRFs) [53, 46]. MRFs organize the nodes/edges as an interconnected system and interpolating between two MRFs captures the joint evolution of the whole graph system. Extending [22], we derive a closed-form Wasserstein distance between graph distributions and leverage it to construct the Bures-Wasserstein (BW) interpolation of two graphs that ensures the OT displacement in graph generation compared to linear interpolation. We then integrate these insights into a flow-matching framework³ called Bures-Wasserstein Flow (BWFlow). Specifically, by defining a probability path via BW interpolation, we obtain smooth, globally coherent velocity fields at intermediate steps (see Fig. 1c) that respect the non-Euclidean, interconnected structure of graphs. Crucially, BWFlow admits simulation-free computation of densities and velocities along the entire path, which translates into efficient, stable training and sampling.

Contributions. **First**, we theoretically and empirically show that the linear interpolation used in existing graph generation models gives suboptimal probability path construction and velocity estimation. **Second**, through parameterizing graphs as MRFs, we introduce BWFlow, a flow-matching model for graph generation that constructs probability paths respecting the graph geometry and develops smooth velocities. **Third**, BWFlow was tested on plain graph and 2D/3D molecule

²Informally, velocity can be viewed as the path tangent. Section 2 gives a more accurate definition.

³For conciseness, we focus on flow models, but the idea is generalizable to diffusions as in Appendix F.2.

generation, demonstrating competitive performance without an excessive search for path manipulation techniques. We further show that BW interpolation consistently outperforms other interpolation methods in building flow matching models, leading to more stable training and sampling convergence.

2 Preliminaries

2.1 Flow matching for graph generation

Flow matching (FM). Generative modeling considers fitting a mapping from state space $\mathcal{S} \rightarrow \mathcal{S}$ that transforms the samples from source distribution, $X_0 \sim p_0$, to samples from target data distribution, $X_1 \sim p_1$ ⁴. Continuous normalizing flow (CNF) [12] parameterizes the transformation through a push-forward equation that interpolates between p_0 and p_1 and constructs a probability path $p_t(\mathcal{X}) = [\psi_t p_0](\mathcal{X})$ through a time-dependent function ψ_t (a.k.a flow). A vector field u_t , defined as $\frac{d}{dt}\psi_t(\mathcal{X}) = u_t(\psi_t(\mathcal{X}))$ with $\psi_0(\mathcal{X}) = \mathcal{X}$, is said to generate p_t if ψ_t satisfies $X_t := \psi_t(X_0) \sim p_t$ for $X_0 \sim p_0$. The FM [33] is designed to match the real velocity field through the loss:

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t, X_t \sim p_t(\cdot)} \|v_\theta(X_t) - u_t(X_t)\|^2. \quad (1)$$

where $v_\theta(\cdot) : \mathcal{S} \rightarrow \mathcal{S}$ is the parameterized velocity field and $t \sim \mathcal{U}[0, 1]$.

Conditional flow matching (CFM). Given that the actual velocity field and the path are not tractable [54], one can construct the per-sample conditional flow. We condition the probability paths on variable $Z \sim \pi(\cdot)$ (for instance, a pair of source and target points $Z = (X_0, X_1)$) and re-write $p_t(\mathcal{X}) = \mathbb{E}_{\pi(\cdot)} p_t(\mathcal{X} | Z)$ and $u_t(\mathcal{X}) = \mathbb{E}_{\pi(\cdot)} u_t(\mathcal{X} | Z)$ where the conditional path and the velocity field are tractable. The CFM aims at regressing a velocity $v_\theta(\cdot)$ to $u_t(\mathcal{X} | Z)$ by the loss,

$$\mathcal{L}_{\text{CFM}}(\theta) := \mathbb{E}_{t, Z \sim \pi(\cdot), p_t(\cdot|Z)} \|v_\theta(X_t) - u_t(X_t | Z)\|^2, \quad (2)$$

where it is shown that the CFM optimization has the same optimum as the FM objective [54].

Graphs as statistical objects. When considering graph generation with CFM, the very first step is to model graphs as statistical objects. For notation, we let $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, \mathcal{X}\}$ denote an undirected graph random variable with edges $\mathcal{E} = \{e_{uv}\}$, nodes $\mathcal{V} = \{v\}$, and node features $\mathcal{X} = \{x_v\}$. A graph realization is denoted as $G = \{V, E, X\} \sim p(\mathcal{G})$. We consider a group of latent variables that controls the graph distribution, specifically the node feature mean $\mathbf{X} = [x_1, x_2, \dots, x_{|\mathcal{V}|}]^\top \in \mathbb{R}^{|\mathcal{V}| \times K}$, the weighted adjacency matrix $\mathbf{W} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$, and the Laplacian matrix $\mathbf{L} = \mathbf{D} - \mathbf{W} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$, with $\mathbf{D} = \text{diag}(\mathbf{W}\mathbf{1})$ being the degree matrix (and $\mathbf{1}$ the all-one vector). In a nutshell, graphs are sampled from $G \sim p(\mathcal{G}; \mathbf{G}) = p(\mathcal{X}, \mathcal{E}; \mathbf{X}, \mathbf{W})$.

Graph generation with CFM. The CFM samples new graphs through iteratively building $G_{t+dt} = G_t + v_t^\theta(G_t) \cdot dt$ with initial $G_0 \sim p_0(\mathcal{G})$ and a trained velocity field $v_t^\theta(G_t)$, so that the medium points follows $G_t \sim p_t(\mathcal{G})$ and terminates at $p_1(\mathcal{G})$. We can parameterize $v_t^\theta(G_t)$ as in [19],

$$v_t^\theta(G_t) = \mathbb{E}_{G_0 \sim p_0(\mathcal{G}), G_1 \sim p_{1|t}^\theta(\cdot|G_t)} [v_t(G_t | G_0, G_1)] \quad (3)$$

As such, training the velocity fields is replaced by training a denoiser $p_{1|t}^\theta(\cdot | G_t)$ to predict the clean datapoint, which is equivalent to maximizing the log-likelihood [45, 9],

$$\mathcal{L}_{\text{CFM}} = \mathbb{E}_{G_1 \sim p_1(\cdot), G_0 \sim p_0(\cdot), t \sim \mathcal{U}_{[0,1]}(\cdot), G_t \sim p_{t|0,1}(\cdot|G_1, G_0)} [\log p_{1|t}^\theta(G_1 | G_t)] \quad (4)$$

where t is sampled from a uniform distribution $\mathcal{U}_{[0,1]}$ and $G_t \sim p_{t|0,1}$ can be obtained in a simulation-free manner. This framework avoids the evaluation of the conditional vector field at training time, which both increases the model robustness and training efficiency.

To proceed, a closed form of $p_t(\cdot | G_0, G_1)$ is required to construct both the probability path and the velocity field $v_t(G_t | G_0, G_1)$. A common selection to decompose the probability density assumes independency for each node and edge [25, 45, 18] giving $p(\mathcal{G}) = p(\mathcal{X})p(\mathcal{E}) = \prod_{v \in \mathcal{V}} p(x_v) \prod_{e_{uv} \in \mathcal{E}} p(e_{uv})$.

⁴For clarity, we denote the calligraphic style \mathcal{X} being the random variable, the plain X the relevant realizations and the bold symbol \mathbf{X} the latent variables (parameters) that controls the distributions, i.e. $X \sim p(\mathcal{X}; \mathbf{X})$.

Choosing $\pi(\cdot) = p_0(\mathcal{G})p_1(\mathcal{G})$, the boundary conditions follow $p_i(\mathcal{G}) = \delta(\mathcal{X}_i = \mathbf{X}_i) \cdot \delta(\mathcal{E}_i = \mathbf{W}_i)$, $\forall i = \{0, 1\}$ with δ the dirac function. This decomposition is further combined with linear interpolation to build the path, as introduced in [54], where,

$$\begin{aligned} p_t(x_v | G_0, G_1) &= \mathcal{N}(t[X_1]_v + (1-t)[X_0]_v, \sigma_t^2) \text{ and } u_t(x_v | X_0, X_1) = [X_1]_v - [X_0]_v, \\ p_t(e_{uv} | G_0, G_1) &= \mathcal{N}(t[E_1]_{uv} + (1-t)[E_0]_{uv}, \sigma_t^2) \text{ and } u_t(e_{uv} | G_0, G_1) = [E_1]_{uv} - [E_0]_{uv}. \end{aligned} \quad (5)$$

Similarly, discrete flow matching frameworks for graph generation [45, 48, 60] is also based on linear interpolation, where the interpolant is sampled from a categorical distribution whose probabilities are simply linear interpolation between the boundary conditions.

2.2 Optimal transport and flow matching

Why linear interpolation? Existing literature [36, 1] argues that the probability path $p_t(\mathcal{X} | Z)$ should be chosen to recover the optimal transport (OT) displacement interpolant [39]. The (Kantorovich) optimal transport problem is to find the transport plan between two probability measures, η_0 and η_1 , with the smallest associated transportation cost defined as follows.

Definition 1 (Wasserstein Distance). Denote the possible coupling as $\pi \in \Pi(\eta_0, \eta_1)$, which is a joint measure on $\mathcal{S} \times \mathcal{S}$ whose marginals are η_0 and η_1 respectively. With $c(X, Y)$ being the cost of transporting the mass between X and Y , the Wasserstein distance is defined as,

$$\mathcal{W}_c(\eta_0, \eta_1) = \inf_{\pi \in \Pi(\eta_0, \eta_1)} \int_{\mathcal{S} \times \mathcal{S}} c(X, Y) d\pi(X, Y). \quad (6)$$

When the data follow Euclidean geometry and both boundary distributions, p_0 and p_1 are described by the Gaussian family, the probability path shown in Eq. (5) with $\sigma_t \rightarrow 0$ becomes a solution to Eq. (6).

As suggested in the motivations, linearly interpolating in the disjoint space of nodes and edges with Eq. (5) does not guarantee the OT displacement in non-Euclidean and interconnected objects like graphs. To overcome the limitation, we utilize Markov random fields to capture the joint evolution of the graph system, and build an FM model that generates graphs with smooth probability paths and consistent velocity.

3 Methodology

In this paper, we introduce Bures–Wasserstein Flow Matching (BWFlow), a novel graph generation framework that is built upon the OT displacement when modeling graphs with Markov Random Fields (MRFs). We begin by casting graphs in an MRF formulation in Section 3.1. We then derive the BWFlow framework in Section 3.2 by formulating and solving the OT displacement problem on the MRF, thereby yielding the fundamental components, interpolations and velocity fields, for FM-based graph generation. Finally, in Section 3.3, we extend BWFlow to discrete FM regimes, enabling its application across a broad spectrum of graph-generation tasks. A schematic overview of the entire BWFlow is illustrated in Fig. 2.

3.1 Graph Markov random fields

We borrow the idea from MRF as a remedy to modeling the complex system organized by graphs, which intrinsically captures the underlying mechanism that jointly generates the nodes and edges. Mathematically, we assume the joint probability density distribution (PDF) of node features and graph structure as $p(\mathcal{G}; \mathbf{G}) = p(\mathcal{X}, \mathcal{E}; \mathbf{X}, \mathbf{W}) = p(\mathcal{X}; \mathbf{X}, \mathbf{W})p(\mathcal{E}; \mathbf{W})$ where the node features and graph structure are interconnected through latent variables \mathbf{X} and \mathbf{W} . For node features \mathcal{X} , we follow the MRF assumption in Zhu et al. [64] and decompose the density into the node-wise potential $\varphi_1(v)$, $\forall v \in \mathcal{V}$ and pair-wise potential $\varphi_2(u, v)$, $\forall e_{uv} \in \mathcal{E}$:

$$p(\mathcal{X}; \mathbf{X}, \mathbf{W}) \propto \prod_v \underbrace{\exp\{-(\nu + d_v)\|\mathbf{V}x_v - \boldsymbol{\mu}_v\|^2\}}_{\varphi_1(v)} \prod_{u,v} \underbrace{\exp\{w_{uv}[(\mathbf{V}x_u - \boldsymbol{\mu}_u)^\top (\mathbf{V}x_v - \boldsymbol{\mu}_v)]\}}_{\varphi_2(u,v)}, \quad (7)$$

with $\|\cdot\|$ the L_2 norm, $(\cdot)^\dagger$ the pseudo-inverse, \mathbf{V} the transformation matrix modulating the graph feature emission, and $\boldsymbol{\mu}_v$ the node-specific latent variable mean. Eq. (7) can be expressed as a colored

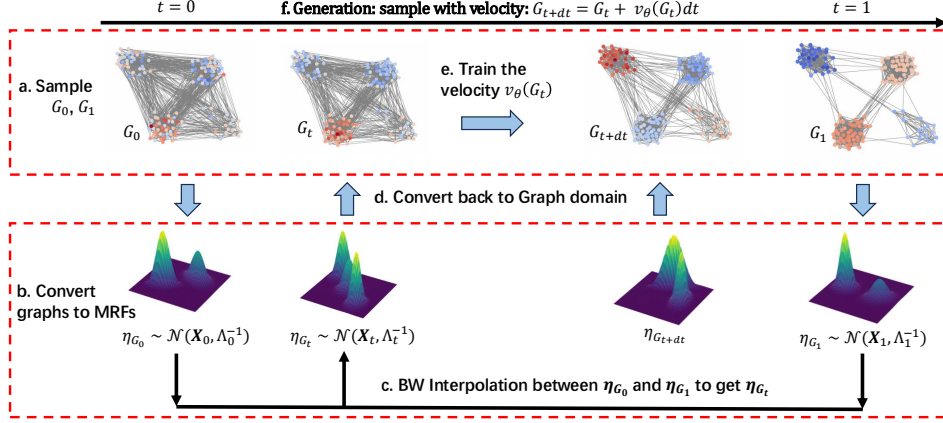


Figure 2: Schematic overview of BWFlow, which consists of: a) Sample the marginal graph condition G_0 and G_1 ; b) Convert graphs to MRFs; c) Interpolate to get intermediate points; d) Convert back to get G_t ; e) Train velocity based on G_t ; and f) Generate new points with the trained velocity.

Gaussian distribution in Eq. (8) given that $Vx_v \sim \mathcal{N}(\mu_v, (\nu I + L)^{-1})$. We further assume that edges are emitted via a Dirac delta, $\mathcal{E} \sim \delta(W)$, yielding our definition of Graph Markov Random Fields (GraphMRF). The derivation can be found in Appendix A.2.

Definition 2 (Graph Markov Random Fields). GraphMRF statistically describes graphs as,

$$p(\mathcal{G}; \mathbf{G}) = p(\mathcal{X}, \mathcal{E}; \mathbf{X}, \mathbf{W}) = p(\mathcal{X}; \mathbf{X}, \mathbf{W}) \cdot p(\mathcal{E}; \mathbf{W}) \text{ where } \mathcal{E} \sim \delta(\mathbf{W}) \text{ and } \text{vec}(\mathcal{X}) \sim \mathcal{N}(\mathbf{X}, \mathbf{\Lambda}^\dagger), \text{ with } \mathbf{X} = \text{vec}(\mathbf{V}^\dagger \boldsymbol{\mu}), \mathbf{\Lambda} = (\nu \mathbf{I} + \mathbf{L}) \otimes \mathbf{V}^\top \mathbf{V}. \quad (8)$$

The \otimes is the Kronecker product, $\text{vec}(\cdot)$ is the vectorization operator and \mathbf{I} is the identity matrix.

Remark 1. GraphMRF explicitly captures node–edge dependencies and preserves the advantages of colored Gaussian distributions. Section 3.2 will soon show that this yields closed-form interpolation and velocity, and the probability path constructed from GraphMRFs remains on the graph manifold that respects the underlying non-Euclidean geometry.

Remark 2. While we emphasize that GraphMRF is not a universal model and imposes certain constraints (Appendix A.3 discusses about the usage scope), it nonetheless captures the dynamics of most graph-generation tasks such as planar, stochastic block models, and molecular graphs.

3.2 Bures-Wasserstein flow matching for graph generation

The optimal transport displacement between graph distributions. Given that the joint probability of graphs decomposed as $p(\mathcal{G}) = p(\mathcal{X}; \mathbf{X}, \mathbf{W})p(\mathcal{E}; \mathbf{W})$ and the measure factorized to $\eta_{\mathcal{G}_j} = \eta_{\mathcal{X}_j} \cdot \eta_{\mathcal{E}_j}$ with $j \in \{0, 1\}$, the graph Wasserstein distance between $\eta_{\mathcal{G}_0}$ and $\eta_{\mathcal{G}_1}$ is written as,

$$(\text{Graph Wasserstein Distance}) \quad d_{\text{BW}}(\mathcal{G}_0, \mathcal{G}_1) := \mathcal{W}_c(\eta_{\mathcal{G}_0}, \eta_{\mathcal{G}_1}) = \mathcal{W}_c(\eta_{\mathcal{X}_0}, \eta_{\mathcal{X}_1}) + \mathcal{W}_c(\eta_{\mathcal{E}_0}, \eta_{\mathcal{E}_1}).$$

We extend Haasler and Frossard [22] and analytically derive the graph Wasserstein distance using the OT formula between Gaussians [16, 42, 52] (see Lemma 2 proved in Appendix B.1) as follows.

Proposition 1 (Bures-Wasserstein Distance). Consider two same-sized graphs $\mathcal{G}_0 \sim p(\mathcal{X}_0, \mathcal{E}_0)$ and $\mathcal{G}_1 \sim p(\mathcal{X}_1, \mathcal{E}_1)$ with \mathbf{V} shared for two graphs, described by the distribution in Definition 2. When the graphs are equipped with graph Laplacian matrices \mathbf{L}_0 and \mathbf{L}_1 satisfying 1) is Positive Semi-Definite (PSD) and 2) has only one zero eigenvalue. The Bures-Wasserstein distance between these two random graph distributions is given by

$$d_{\text{BW}}(\mathcal{G}_0, \mathcal{G}_1) = \|\mathbf{X}_0 - \mathbf{X}_1\|_F^2 + \beta \text{trace} \left(\mathbf{L}_0^\dagger + \mathbf{L}_1^\dagger - 2 \left(\mathbf{L}_0^{\dagger/2} \mathbf{L}_1^\dagger \mathbf{L}_0^{\dagger/2} \right)^{1/2} \right), \quad (9)$$

as $\nu \rightarrow 0$ and β is a constant related to the norm of \mathbf{V}^\dagger . The proof can be found in Appendix B.2.

Based on the Bures-Wasserstein (BW) distance, we then derive the OT interpolant for two graphs, which is the solution of the displacement minimization problem described as,

$$\mathcal{G}_t = \arg \min_{\tilde{\mathcal{G}}} (1-t)d_{\text{BW}}(\mathcal{G}_0, \tilde{\mathcal{G}}) + td_{\text{BW}}(\tilde{\mathcal{G}}, \mathcal{G}_1). \quad (10)$$

The probability path. The interpolation is obtained through solving Eq. (10) with the BW distance defined in Proposition 1, we prove the minimizer of the above problem has the form in Proposition 2. The proof can be found in Appendix C.1.

Proposition 2 (Bures-Wasserstein interpolation). *The graph minimizer of Eq. (10), $\mathcal{G}_t = \{\mathcal{V}, \mathcal{E}_t, \mathcal{X}_t\}$, have its node features following a colored Gaussian distribution, $\mathcal{X}_t \sim \mathcal{N}(\mathbf{X}_t, \mathbf{\Lambda}_t^\dagger)$ with $\mathbf{\Lambda}_t = (\nu \mathbf{I} + \mathbf{L}_t) \otimes \mathbf{V}^\top \mathbf{V}$ and edges following $\mathcal{E}_t \sim \delta(\mathbf{W}_t)$, specifically,*

$$\mathbf{L}_t^\dagger = \mathbf{L}_0^{1/2} \left((1-t)\mathbf{L}_0^\dagger + t \left(\mathbf{L}_0^{\dagger/2} \mathbf{L}_1^\dagger \mathbf{L}_0^{\dagger/2} \right)^{1/2} \right)^2 \mathbf{L}_0^{1/2}, \quad \mathbf{X}_t = (1-t)\mathbf{X}_0 + t\mathbf{X}_1 \quad (11)$$

The interpolant provides a closed form for the induced probability path $p(\mathcal{G}_t | G_0, G_1)$ and the velocity $v(\mathcal{G}_t | G_0, G_1)$ that is easy to access without any simulation.

The velocity. We consider the reparameterization as in Eq. (3) and derive the conditional velocity $v_t(\mathcal{G}_t | G_1, G_0)$ as in Proposition 3.

Proposition 3 (Bures-Wasserstein velocity). *For the graph \mathcal{G}_t following BW interpolation in Proposition 2, the conditional velocity at time t with observation \mathcal{G}_t is given as,*

$$v_t(\mathcal{E}_t | G_0, G_1) = \dot{\mathbf{W}}_t = \text{diag}(\dot{\mathbf{L}}_t) - \dot{\mathbf{L}}_t, \quad v_t(\mathbf{X}_t | G_0, G_1) = \frac{1}{1-t}(\mathbf{X}_1 - \mathbf{X}_t) \quad (12)$$

$$\text{with } \dot{\mathbf{L}}_t = 2\mathbf{L}_t - \mathbf{T}\mathbf{L}_t - \mathbf{L}_t\mathbf{T} \text{ and } \mathbf{T} = \mathbf{L}_0^{1/2}(\mathbf{L}_0^{\dagger/2} \mathbf{L}_1^\dagger \mathbf{L}_0^{\dagger/2})^{1/2} \mathbf{L}_0^{1/2}$$

where $\mathbf{W}_t = \mathbf{D}_t - \mathbf{L}_t$ and \mathbf{L}_t defined in Eq. (11). Derivation can be found in Appendix C.2.

With Proposition 2 and Proposition 3, we are now able to formally construct the algorithms for Bures-Wasserstein flow matching. Taking continuous flow matching as an example, Algorithms 1 and 2 respectively introduce the training and sampling pipelines for our BWFlow.

Remark: The BW interpolation and velocity both deviate from the linear flow matching framework and require extra computational cost. However, there exist multiple ways to analytically calculate or numerically approximate the velocity for training and inference. The choice of these methods depends on the trade-off between training stability, sampling efficiency, etc. In Appendix E, we provide a discussion about the design space of BW interpolation and velocity.

3.3 Discrete Bures-Wasserstein flow matching for graph generation

Up to now we are working on the scenario when $p(\mathcal{X} | \mathbf{X}, \mathbf{W})$ is a Gaussian and $p(\mathcal{E} | \mathbf{W})$ is a Dirac distribution. However, previous studies have observed a significant improvement of the discrete counterpart of the continuous graph generation models [55, 60, 45]. To benefit our model from such a nature, we derive the discrete Bures-Wasserstein flow matching for graph generation.

The discrete probability path. We design the probability path as discrete distributions,

$$p_t(x_v | G_0, G_1) = \text{Categorical}([\mathbf{X}_t]_v), \quad p_t(e_{uv} | G_0, G_1) = \text{Bernoulli}([\mathbf{W}_t]_{uv}) \quad (13)$$

$$\text{s.t. } p_0(\mathcal{G}) = \delta(G_0, \cdot), p_1(\mathcal{G}) = \delta(G_1, \cdot)$$

where $\mathbf{W}_t = \mathbf{D}_t - \mathbf{L}_t$ with \mathbf{X}_t and \mathbf{L}_t defined the same in Eq. (11). We consider the fact that the Dirac distribution is a special case when the Categorical/Bernoulli distribution has probability 1 or 0, so the boundary condition $p_0(\mathcal{G}) = \delta(G_0, \cdot), p_1(\mathcal{G}) = \delta(G_1, \cdot)$ holds. Even though we are not sampling from Gaussian distributions anymore, it is possible to approximate the Wasserstein distance between two multivariate discrete distributions with the Gaussian counterpart so the conclusions, such as optimal transport displacements, still hold. We left the discussion in Appendix D.2.

The discrete velocity fields. The path of node features \mathcal{X}_t can be re-written as $p_t(\mathcal{X}) = (1-t)\delta(\cdot, \mathbf{X}_0) + t\delta(\cdot, \mathbf{X}_1)$ so the conditional velocity can be accessed through $v_t(\mathbf{X}_t | G_0, G_1) =$

Algorithm 1: BWFlow Training

Input: Ref. dist p_0 and dataset $\mathcal{D} \sim p_1$.**Output:** Trained model $f_\theta(G_t, t)$.

```
1 Initialize model  $f_\theta(G_t, t)$ ;  
2 while  $f_\theta$  not converged do  
    /* Sample Boundary Graphs */  
3    Sample batched  $\{G_0\} \sim p_0, \{G_1\} \sim \mathcal{D}$ ;  
    /* Prob.path Construction */  
4    Sample  $t \sim \mathcal{U}(0, 1)$ ;  
5    Calculate the BW interpolation  
     $p(G_t | G_0, G_1)$  via Eq. (11);  
    /* Denoising - x-prediction */  
6     $p_{1|t}^\theta(\cdot | G_t) \leftarrow f_\theta(G_t, t)$ ;  
7    Loss calculation via Eq. (4);  
8    optimizer.step();
```

Algorithm 2: BWFlow Sampling

Input: Reference distribution p_0 , TrainedModel $f_\theta(G_t, t)$, Small time step dt ,**Output:** Generated Graphs $\{\hat{G}_1\}$.

```
1 Initialize samples  $\{\hat{G}_0\} \sim p_0$ ;  
2 Initialize the model  $p_{1|t}^\theta(\cdot | G_t) \leftarrow f_\theta(G_t, t)$   
   for  $t \leftarrow 0$  to  $1 - dt$  by  $dt$  do  
   /* Denoising - x-prediction */  
3   Predict  $\tilde{G}_1 \leftarrow p_{1|t}^\theta(\cdot | \hat{G}_t)$ ;  
   /* Velocity calculation */  
4   Calculate  $v_\theta(\hat{G}_t | \hat{G}_0, \tilde{G}_1)$  via Eq. (12);  
   /* Numerical Sampling */  
5   Sample  $\hat{G}_{t+dt} \sim \hat{G}_t + v_\theta(\hat{G}_t)dt$ 
```

$[\delta(\cdot, \mathbf{X}_1) - \delta(\cdot, \mathbf{X}_t)] / (1 - t)$. However, the probability path of edges \mathcal{E}_t , shown in Eqs. (11) and (13), cannot be written as a mixture of two boundary conditions given the non-linear interpolation. To this end, we derive in Appendix D.3 that the discrete velocity follows,

$$v_t(E_t | G_1, G_0) = (1 - 2E_t) \frac{\dot{\mathbf{W}}_t}{\mathbf{W}_t \circ (1 - \mathbf{W}_t)}, \quad (14)$$

where $\mathbf{W}_t = \mathbf{D}_t - \mathbf{L}_t$, $\dot{\mathbf{W}}_t = \text{diag}(\dot{\mathbf{L}}_t) - \dot{\mathbf{L}}_t$ with $\mathbf{L}_t, \dot{\mathbf{L}}_t$ defined in Eqs. (11) and (12) respectively. With the interpolation and velocity defined, the discrete flow matching is built in Algorithms 3 and 4.

4 Experiments

We evaluate the BWFlow algorithms through both the plain graph generation and real-world molecule generation tasks. We first outline the experimental setup in Section 4.1, followed by a general comparison in Section 4.2. Next, we conduct behavior analysis on the impact of interpolation methods and the corresponding velocity construction on graph generation performance in Section 4.3, which demonstrates the effectiveness and benefit of flow along Bures-Wasserstein interpolation.

4.1 Experiment settings

Dataset. For plain graph generation, we evaluate the quality of generated graphs on three benchmark datasets following previous works [37, 55, 4], specifically, **planar** graphs, **tree** graphs, and stochastic blocking models (**SBM**). Two datasets, namely MOSES [43] and GUACAMOL [7], are benchmarked to test the model performance on 2D molecule generation. For 3D molecule generation with coordinate data, we test the model on QM9 [47] and GEOM-DRUGS [2].

Metrics. In plain graph generation, the evaluation metrics include the percentage of Valid, Unique, and Novel (**V.U.N.**) graphs, and the average maximum mean discrepancy ratio (**A.Ratio**) of graph statistics between the set of generated graphs and the test set are reported (details in Appendix I.3). For molecule generation, we test two scenarios with and without bond type information, where the latter validates the capacity of our methods in generating the graph structures. To this end, we develop a new relaxed metric to measure the stability and validity of atoms and molecules when bond types are not available. Specifically, the atom-wise stability is relaxed as:

$$\text{Stability of Atom } i: s_i = \mathbb{I}[\exists \{b_{ij}\}_{j \in \mathcal{N}_i} \in \prod_{j \in \mathcal{N}_i} B_{ij} : \sum_{j \in \mathcal{N}_i} b_{ij} = \text{EV}_i], \text{ with the identity function } \mathbb{I}.$$

This means atom i is “relaxed-stable” if there is at least one way to pick allowed bond types (B_{ij}) to its neighbors \mathcal{N}_i so that their total exactly matches the expected valences EV_i . Such a relaxed stability of atoms (**Atom.Stab.**) inherently defines molecule stability (**Mol.Stab.**) and the **validity** of a molecule, which are the shared metrics for both 2D/3D molecule generation. In addition to these

metrics, distribution metrics are also used for 2D molecules (FCD, Scaf, etc.), and 3D generations (charge distributions, atom total variation, angles, etc.). Details in Appendix I.3.

Setup. To isolate the impact from model architecture, we follow Qin et al. [45] to fix the backbone model as the same graph transformers. It is shown that sampling/training distortion and target guidance have a significant impact on the performance of graph generation tasks [45]. In our experiment, the best model performance is obtained with these technologies, but in behavior analysis, we disabled time distortion and target guidance for a fair comparison. In molecular generation, two scenarios with and without bond types information are considered to better evaluate the ability of generating graph structures. More experimental details can be found in Appendix I.1.

4.2 Main results for graph generation

Plain graph generation. In Table 1, we report both V.U.N. and A.Ratio. As performance on these benchmarks continues to fluctuate significantly even after convergence and the results are near saturated, we present not only the best scores but also the exponentially moving averaged (EMA) results on last 5 checkpoints and decay 0.999. BWFlow outperforms most competitors on Planar and SBM graph generation. The lone exception is the tree graphs, where our model falls short. We attribute this gap to the fundamentally different generation process for tree graphs (which reside in hyperbolic space [61]) and thus violate our MRF assumptions.

Table 1: Plain Graph generation performance. We sampled 5 times (each run generates 40 graphs) to calculate the mean and standard deviation. We only keep the main diffusion/flow model for comparison, while other models are included in the full version at Table 8.

Model	Class	Planar		Tree		SBM	
		V.U.N. \uparrow	A.Ratio \downarrow	V.U.N. \uparrow	A.Ratio \downarrow	V.U.N. \uparrow	A.Ratio \downarrow
Train set	—	100	1.0	100	1.0	85.9	1.0
DiGress [55]	Diffusion	77.5	5.1	90.0	1.6	60.0	1.7
HSpectre [4]	Diffusion	95.0	2.1	100.0	4.0	75.0	10.5
GruM [28]	Diffusion	90.0	1.8	—	—	85.0	1.1
CatFlow [18]	Flow	80.0	—	—	—	85.0	—
DeFoG [45]	Flow	99.5 \pm 1.0	1.6 \pm 0.4	96.5 \pm 2.6	1.6 \pm 0.4	90.0 \pm 5.1	4.9 \pm 1.3
BWFlow	Flow	97.5 \pm 2.5	1.3 \pm 0.4	75.5 \pm 2.4	1.4 \pm 0.3	90.5 \pm 4.0	3.8 \pm 0.9
DeFoG (EMA)	Flow	77.5 \pm 8.37	3.5 \pm 1.7	73.1 \pm 11.4	1.50 \pm 0.3	85.0 \pm 7.1	3.7 \pm 0.9
BWFlow (EMA)	Flow	84.8 \pm 6.44	2.4 \pm 0.9	68.1 \pm 12.2	1.24 \pm 0.2	83.5 \pm 6.0	2.4 \pm 0.6

2D molecule graph generation. The model performance is illustrated in Table 2. In both datasets, BWFlow can achieve competitive results near the state-of-the-art (SOTA) flow matching models [45], and outperforms the diffusion models. Given that MOSES and GUACAMOL benchmarks are approaching saturation, the fact that BWFlow achieves results on par with the SOTA models serves as strong evidence of its effectiveness.

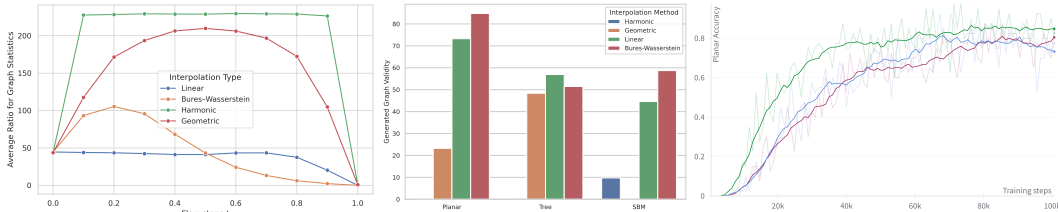
Table 2: Large molecule generation results. Only diffusion and flow models are reported. Table 15 gives further experiments with binary edge types.

Model	Guacamol				MOSES						
	Val. \uparrow	V.U. \uparrow	V.U.N. \uparrow	FCD \uparrow	Val. \uparrow	Unique. \uparrow	Novelty \uparrow	Filters \uparrow	FCD \downarrow	SNN \uparrow	Scaf \uparrow
Training set	100.0	100.0	0.0	92.8	100.0	100.0	0.0	100.0	0.01	0.64	99.1
DiGress [55]	85.2	85.2	85.1	68.0	85.7	100.0	95.0	97.1	1.19	0.52	14.8
DisCo [60]	86.6	86.6	86.5	59.7	88.3	100.0	97.7	95.6	1.44	0.50	15.1
Cometh [48]	98.9	98.9	97.6	72.7	90.5	99.9	92.6	99.1	1.27	0.54	16.0
DeFoG [45]	99.0	99.0	97.9	73.8	92.8	99.9	92.1	98.9	1.95	0.55	14.4
BWFlow (Ours)	98.8	98.9	97.4	69.2	92.0	100.0	94.5	98.4	1.32	0.56	15.3

3D molecule generation. Table 3 gives the results on the 3D molecule generation task with explicit hydrogen, where we ignore the bond type but just view the adjacency matrix as a binary one for validating the power of generating graph structures. Interestingly, the empirical results show that even without edge type, the 3D graph generation model already can capture the molecule data distribution. And our BWFlow significantly outperforms the SOTA models, including MiDi [56] and FlowMol [17]. We believe a promising future direction is to incorporate the processing of multiple bond types into our framework, which would potentially raise the performance by a margin.

Dataset	Interpolation	Metrics							
		μ	V.U.N(%)	Mol.Stab.	Atom.Stab.	Connected(%)	Charge(10^{-2})	Atom(10^{-2})	Angles($^{\circ}$)
QM9 (with h)	MiDi	1.01	93.13	93.98	99.60	99.21	0.2	3.7	2.21
	FlowMol	1.01	87.53	88.45	99.13	99.09	0.4	4.2	2.72
	BWFlow	1.01	96.45	97.84	99.84	99.24	0.1	2.3	1.96
GEOM (with h)	Midi	1.34	78.23	32.42	89.61	79.15	0.6	11.2	9.6
	FlowMol	1.34	82.20	36.90	94.60	59.98	0.4	8.8	6.5
	BWFlow	1.20	87.75	46.80	95.08	73.53	0.1	6.5	3.96

Table 3: Quantitative experimental results on 3D Molecule Generation with explicit hydrogen.



(a) The evolution of graph statistics ratio along the probability path. (b) The impact of interpolation methods on the performance. (c) Convergence analysis of BW-Flow and flows with linear interpolations.

Figure 3: Ablation studies for Bures-Wasserstein Flow Matching.

4.3 Behavior analysis

BWFlow provides smooth velocity in probability paths. To illustrate how BWFlow models the smooth evolution of graphs, we compute the A.Ratio on SBM datasets (the figures for the others are in Fig. 5) between generated graph interpolants and test data for $t \in [0, 1]$, as shown in Fig. 3a. In contrast to the linear (arithmetic) interpolation, BW interpolation initially exposes the model to more out-of-distribution samples with increased A.Ratio. After this early exploration, the A.Ratio monotonously converges, yielding a smooth interpolation between the reference graphs and the data points. This behavior enhances both the model robustness and velocity estimation, which helps in covering the convergence gap in the generation stage as in Fig. 1c. In comparison, harmonic and geometric interpolations step outside the valid graph domain, making the learning ill-posed.

The impact of interpolation metrics on the model performance. Fig. 3b illustrated a bar plot that compares interpolation methods on the ability of generating valid plain graphs measured by V.U.N., which shows the superiority of BW interpolation in capturing graph distributions. Fig. 3c illustrated an example (in planar graph generation) of the convergence curve at the training stage (full results in Table 10), which suggests that BWFlow can bring a faster convergence speed compared to FM methods constructed with linear (arithmetic) interpolations.

5 Discussion and future work

In this paper, we introduce BWFlow, a flow matching model that captures the non-Euclidean and interconnected properties of graphs. While we show BWFlow exhibits outstanding performance in various graph generation tasks, it faces the following limitations that motivate solid future work.

Extension to multiple relation types. As our framework is built upon the interpolation parameterized by the Graph Laplacian, it is not easily generalizable to the graph generation with multiple edge types. We made preliminary attempts at the extension but a comprehensive design is still required.

Lower computational complexity. While constructing the probability path and the velocity, our BW interpolation suffers from an extra cost due to its request to compute the pseudo-inverse of the Laplacian. Compared to linear interpolation, this brings $O(N^3)$ extra complexity theoretically, and empirically 2x training time and inference time, which is non-negligible in large graph generation. We aim to develop iterative optimization methods to make training faster in future work.

More universal interpolation that accommodates the geometry. In our experiments on the tree dataset, performance was unsatisfactory. We attribute this issue to the unique geometry of tree-structured graphs. A promising future work includes selecting the adaptive interpolation schemes that accommodate the intrinsic geometry of a graph.

References

- [1] Michael S. Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. In *ICLR*. OpenReview.net, 2023.
- [2] Simon Axelrod and Rafael Gómez-Bombarelli. GEOM: energy-annotated molecular conformations for property prediction and molecular generation. *CoRR*, abs/2006.05531, 2020.
- [3] Eric Bach, Simon Rogers, John Williamson, and Juho Rousu. Probabilistic framework for integration of mass spectrum and retention time information in small molecule identification. *Bioinformatics*, 37(12):1724–1731, 11 2020. ISSN 1367-4803. doi: 10.1093/bioinformatics/btaa998. URL <https://doi.org/10.1093/bioinformatics/btaa998>.
- [4] Andreas Bergmeister, Karolis Martinkus, Nathanaël Perraudin, and Roger Wattenhofer. Efficient and scalable graph generation through iterative local expansion. In *ICLR*. OpenReview.net, 2024.
- [5] Rajendra Bhatia, Tanvi Jain, and Yongdo Lim. On the bures–wasserstein distance between positive definite matrices. *Expositiones Mathematicae*, 37(2):165–191, 2019. ISSN 0723-0869. doi: <https://doi.org/10.1016/j.exmath.2018.01.002>. URL <https://www.sciencedirect.com/science/article/pii/S0723086918300021>.
- [6] Camille Bilodeau, Wengong Jin, Tommi Jaakkola, Regina Barzilay, and Klavs F Jensen. Generative models for molecular discovery: Recent advances and challenges. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 12(5):e1608, 2022.
- [7] Nathan Brown, Marco Fiscato, Marwin H. S. Segler, and Alain C. Vaucher. Guacamol: Benchmarking models for de novo molecular design. *J. Chem. Inf. Model.*, 59(3):1096–1108, 2019.
- [8] Andrew Campbell, Joe Benton, Valentin De Bortoli, Thomas Rainforth, George Deligiannidis, and Arnaud Doucet. A continuous time framework for discrete denoising models. In *NeurIPS*, 2022.
- [9] Andrew Campbell, Jason Yim, Regina Barzilay, Tom Rainforth, and Tommi S. Jaakkola. Generative flows on discrete state-spaces: Enabling multimodal flows with applications to protein co-design. In *ICML*. OpenReview.net, 2024.
- [10] Nicola De Cao and Thomas Kipf. Molgan: An implicit generative model for small molecular graphs. *CoRR*, abs/1805.11973, 2018.
- [11] Ricky T. Q. Chen and Yaron Lipman. Flow matching on general geometries. In *ICLR*. OpenReview.net, 2024.
- [12] Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. In *NeurIPS*, pages 6572–6583, 2018.
- [13] Xiaohui Chen, Jiaying He, Xu Han, and Liping Liu. Efficient and degree-guided graph generation via discrete diffusion modeling. In *ICML*, volume 202 of *Proceedings of Machine Learning Research*, pages 4585–4610. PMLR, 2023.
- [14] Hanjun Dai, Azade Nazi, Yujia Li, Bo Dai, and Dale Schuurmans. Scalable deep generative modeling for sparse graphs. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 2302–2312. PMLR, 2020.
- [15] Nathaniel Lee Diamant, Alex M. Tseng, Kangway V. Chuang, Tommaso Biancalani, and Gabriele Scalia. Improving graph generation by restricting graph bandwidth. In *ICML*, volume 202 of *Proceedings of Machine Learning Research*, pages 7939–7959. PMLR, 2023.
- [16] D.C Dowson and B.V Landau. The fréchet distance between multivariate normal distributions. *Journal of Multivariate Analysis*, 12(3):450–455, 1982. ISSN 0047-259X. doi: [https://doi.org/10.1016/0047-259X\(82\)90077-X](https://doi.org/10.1016/0047-259X(82)90077-X). URL <https://www.sciencedirect.com/science/article/pii/0047259X8290077X>.
- [17] Ian Dunn and David Ryan Koes. Mixed continuous and categorical flow matching for 3d de novo molecule generation. *ArXiv*, pages arXiv–2404, 2024.

- [18] Floor Eijkelboom, Grigory Bartosh, Christian Andersson Naesseth, Max Welling, and Jan-Willem van de Meent. Variational flow matching for graph generation. In *NeurIPS*, 2024.
- [19] Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky T. Q. Chen, Gabriel Synnaeve, Yossi Adi, and Yaron Lipman. Discrete flow matching. In *NeurIPS*, 2024.
- [20] Nikhil Goyal, Harsh Vardhan Jain, and Sayan Ranu. Graphgen: A scalable approach to domain-agnostic labeled graph generation. In *WWW*, pages 1253–1263. ACM / IW3C2, 2020.
- [21] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *KDD*, pages 855–864. ACM, 2016.
- [22] Isabel Haasler and Pascal Frossard. Bures-wasserstein means of graphs. In *AISTATS*, volume 238 of *Proceedings of Machine Learning Research*, pages 1873–1881. PMLR, 2024.
- [23] Kilian Konstantin Haefeli, Karolis Martinkus, Nathanaël Perraudin, and Roger Wattenhofer. Diffusion models for graphs benefit from discrete state spaces. In *The First Learning on Graphs Conference*, 2022. URL <https://openreview.net/forum?id=CtsKBwhTMKg>.
- [24] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.
- [25] Xiaoyang Hou, Tian Zhu, Milong Ren, Dongbo Bu, Xin Gao, Chunming Zhang, and Shiwei Sun. Improving molecular graph generation with flow matching and optimal transport. *CoRR*, abs/2411.05676, 2024.
- [26] John Ingraham, Vikas K. Garg, Regina Barzilay, and Tommi S. Jaakkola. Generative models for graph-based protein design. In *NeurIPS*, pages 15794–15805, 2019.
- [27] Keyue Jiang, Bohan Tang, Xiaowen Dong, and Laura Toni. Heterogeneous graph structure learning through the lens of data-generating processes. In *The 28th International Conference on Artificial Intelligence and Statistics*, 2025. URL <https://openreview.net/forum?id=JHK0QBkdYY>.
- [28] Jaehyeong Jo, Dongki Kim, and Sung Ju Hwang. Graph generation with diffusion mixture. In *ICML*. OpenReview.net, 2024.
- [29] Kacper Kapusniak, Peter Potaptchik, Teodora Reu, Leo Zhang, Alexander Tong, Michael M. Bronstein, Avishek Joey Bose, and Francesco Di Giovanni. Metric flow matching for smooth interpolations on the data manifold. In *NeurIPS*, 2024.
- [30] Thomas N. Kipf and Max Welling. Variational graph auto-encoders. *CoRR*, abs/1611.07308, 2016.
- [31] Mufei Li, Eleonora Kreacic, Vamsi K. Potluru, and Pan Li. Graphmaker: Can diffusion models generate large attributed graphs? *CoRR*, abs/2310.13833, 2023.
- [32] Renjie Liao, Yujia Li, Yang Song, Shenlong Wang, William L. Hamilton, David Duvenaud, Raquel Urtasun, and Richard S. Zemel. Efficient graph generation with graph recurrent attention networks. In *NeurIPS*, pages 4257–4267, 2019.
- [33] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=PqvMRDCJT9t>.
- [34] Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky T. Q. Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code. *CoRR*, abs/2412.06264, 2024.
- [35] Chengyi Liu, Wenqi Fan, Yunqing Liu, Jiatong Li, Hang Li, Hui Liu, Jiliang Tang, and Qing Li. Generative diffusion models on graphs: Methods and applications. In *IJCAI*, pages 6702–6711. ijcai.org, 2023.
- [36] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *ICLR*. OpenReview.net, 2023.

- [37] Karolis Martinkus, Andreas Loukas, Nathanaël Perraudin, and Roger Wattenhofer. SPECTRE: spectral conditioning helps to overcome the expressivity limits of one-shot graph generators. In *ICML*, volume 162 of *Proceedings of Machine Learning Research*, pages 15159–15179. PMLR, 2022.
- [38] Karolis Martinkus, Andreas Loukas, Nathanaël Perraudin, and Roger Wattenhofer. SPECTRE: spectral conditioning helps to overcome the expressivity limits of one-shot graph generators. In *ICML*, volume 162 of *Proceedings of Machine Learning Research*, pages 15159–15179. PMLR, 2022.
- [39] Robert J. McCann. A convexity principle for interacting gases. *Advances in Mathematics*, 128(1):153–179, 1997. ISSN 0001-8708. doi: <https://doi.org/10.1006/aima.1997.1634>. URL <https://www.sciencedirect.com/science/article/pii/S0001870897916340>.
- [40] Giorgia Minello, Alessandro Bicciato, Luca Rossi, Andrea Torsello, and Luca Cosmo. Generating graphs via spectral diffusion. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=AAxBfJNHdt>.
- [41] Chenhao Niu, Yang Song, Jiaming Song, Shengjia Zhao, Aditya Grover, and Stefano Ermon. Permutation invariant graph generation via score-based generative modeling. In *AISTATS*, volume 108 of *Proceedings of Machine Learning Research*, pages 4474–4484. PMLR, 2020.
- [42] I. Olkin and F. Pukelsheim. The distance between two random vectors with given dispersion matrices. *Linear Algebra and its Applications*, 48:257–263, 1982. ISSN 0024-3795. doi: [https://doi.org/10.1016/0024-3795\(82\)90112-4](https://doi.org/10.1016/0024-3795(82)90112-4). URL <https://www.sciencedirect.com/science/article/pii/0024379582901124>.
- [43] Daniil Polykovskiy, Alexander Zhebrak, Benjamín Sánchez-Lengeling, Sergey Golovanov, Oktai Tatanov, Stanislav Belyaev, Rauf Kurbanov, Aleksey Artamonov, Vladimir Aladinskiy, Mark Veselov, Artur Kadurin, Sergey I. Nikolenko, Alán Aspuru-Guzik, and Alex Zhavoronkov. Molecular sets (MOSES): A benchmarking platform for molecular generation models. *CoRR*, abs/1811.12823, 2018.
- [44] Aram-Alexandre Pooladian, Heli Ben-Hamu, Carles Domingo-Enrich, Brandon Amos, Yaron Lipman, and Ricky T. Q. Chen. Multisample flow matching: Straightening flows with minibatch couplings. In *ICML*, volume 202 of *Proceedings of Machine Learning Research*, pages 28100–28127. PMLR, 2023.
- [45] Yiming Qin, Manuel Madeira, Dorina Thanou, and Pascal Frossard. Defog: Discrete flow matching for graph generation. *CoRR*, abs/2410.04263, 2024.
- [46] Meng Qu, Yoshua Bengio, and Jian Tang. GMNN: graph markov neural networks. In *ICML*, volume 97 of *Proceedings of Machine Learning Research*, pages 5241–5250. PMLR, 2019.
- [47] Raghunathan Ramakrishnan, Pavlo O Dral, Matthias Rupp, and O Anatole Von Lilienfeld. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific data*, 1(1):1–7, 2014.
- [48] Antoine Siraudin, Fragkiskos D. Malliaros, and Christopher Morris. Cometh: A continuous-time discrete-state graph diffusion model, 2024. URL <https://arxiv.org/abs/2406.06449>.
- [49] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*. OpenReview.net, 2021.
- [50] Hannes Stärk, Bowen Jing, Chenyu Wang, Gabriele Corso, Bonnie Berger, Regina Barzilay, and Tommi S. Jaakkola. Dirichlet flow matching with applications to DNA sequence design. In *ICML*. OpenReview.net, 2024.
- [51] Haoran Sun, Lijun Yu, Bo Dai, Dale Schuurmans, and Hanjun Dai. Score-based continuous-time discrete diffusion models. In *ICLR*. OpenReview.net, 2023.
- [52] Asuka Takatsu. On wasserstein geometry of gaussian measures. *Probabilistic approach to geometry*, 57:463–472, 2010.

- [53] Ben Taskar, Pieter Abbeel, Ming-Fai Wong, and Daphne Koller. Relational markov networks. *Introduction to statistical relational learning*, 175:200, 2007.
- [54] Alexander Tong, Kilian Fatras, Nikolay Malkin, Guillaume Huguët, Yanlei Zhang, Jarrid Rector-Brooks, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models with minibatch optimal transport. *Trans. Mach. Learn. Res.*, 2024, 2024.
- [55] Clément Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pascal Frossard. Digress: Discrete denoising diffusion for graph generation. In *ICLR*. OpenReview.net, 2023.
- [56] Clément Vignac, Nagham Osman, Laura Toni, and Pascal Frossard. Midi: Mixed graph and 3d denoising diffusion for molecule generation. In *ECML/PKDD (2)*, volume 14170 of *Lecture Notes in Computer Science*, pages 560–576. Springer, 2023.
- [57] C. Villani and American Mathematical Society. *Topics in Optimal Transportation*. Graduate studies in mathematics. American Mathematical Society, 2003. ISBN 9781470418045. URL <https://books.google.co.uk/books?id=MyPjjgEACAAJ>.
- [58] Fu-Yun Wang, Ling Yang, Zhaoyang Huang, Mengdi Wang, and Hongsheng Li. Rectified diffusion: Straightness is not your need in rectified flow. *CoRR*, abs/2410.07303, 2024.
- [59] Martin Weigt, Robert A. White, Hendrik Szurmant, James A. Hoch, and Terence Hwa. Identification of direct residue contacts in protein–protein interaction by message passing. *Proceedings of the National Academy of Sciences*, 106(1):67–72, 2009. doi: 10.1073/pnas.0805923106. URL <https://www.pnas.org/doi/abs/10.1073/pnas.0805923106>.
- [60] Zhe Xu, Ruizhong Qiu, Yuzhong Chen, Huiyuan Chen, Xiran Fan, Menghai Pan, Zhichen Zeng, Mahashweta Das, and Hanghang Tong. Discrete-state continuous-time diffusion for graph generation. *arXiv preprint arXiv:2405.11416*, 2024.
- [61] Menglin Yang, Min Zhou, Zhihao Li, Jiahong Liu, Lujia Pan, Hui Xiong, and Irwin King. Hyperbolic graph neural networks: A review of methods and applications. *CoRR*, abs/2202.13852, 2022.
- [62] Jiaxuan You, Rex Ying, Xiang Ren, William L. Hamilton, and Jure Leskovec. Graphrnn: Generating realistic graphs with deep auto-regressive models. In *ICML*, volume 80 of *Proceedings of Machine Learning Research*, pages 5694–5703. PMLR, 2018.
- [63] Meng Yu and Kun Zhan. Bias mitigation in graph diffusion models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=CSj72Rr2PB>.
- [64] Xiaojin Zhu, John Lafferty, and Zoubin Ghahramani. *Semi-supervised learning: From Gaussian fields to Gaussian processes*. School of Computer Science, Carnegie Mellon University, 2003.
- [65] Yanqiao Zhu, Yuanqi Du, Yinkai Wang, Yichen Xu, Jieyu Zhang, Qiang Liu, and Shu Wu. A survey on deep graph generation: Methods and applications. In *LoG*, volume 198 of *Proceedings of Machine Learning Research*, page 47. PMLR, 2022.

A Graph Markov random fields: background and theory

A.1 Background of Markov random fields

Markov random fields (MRFs) were originally developed to describe the dynamics of interconnected physical systems such as molecules and proteins [59, 3]. MRFs are energy-based models that have the following probability density:

$$p(\mathcal{X}) = \frac{1}{Z} \prod_c \phi_c(x_c) = \frac{1}{Z} e^{-E(\mathcal{X})/kT}, \quad (15)$$

where the energy $E(\mathcal{X})$ is used to describe the whole connected system. For instance, in the molecule system that consists of atoms and bonds, the overall energy is decomposed into the atom-wise potential $\varphi_1(v)$, $\forall v \in \mathcal{V}$ and bond-wise potential $\varphi_2(u, v)$, $\forall e_{uv} \in \mathcal{E}$. MRFs serve as a natural and elegant way to describe general graph systems.

The energy-based models have an intrinsic relationship with generative models. As an example, Song et al. [49] derived the relationship between diffusion models and the Langevin dynamics, which is used to describe the evolution of an energy-based model. It is shown that the diffusion models are trying to approximate the score function $\nabla_{\mathcal{X}} \log p(\mathcal{X})$. In the energy-based models, the score function is just the gradient of energy, $\nabla_{\mathcal{X}} \log p(\mathcal{X}) = -\nabla_{\mathcal{X}} E(\mathcal{X})$, and the Langevin dynamics becomes transiting between states with different energies.

The idea of our paper originated from the two facts: MRFs are energy-based model describing connected systems, and the energy-based models have their intrinsic relationship with the diffusion/flow models. Thus, if a model is required to describe the evolution of the whole graph system, we believe it is natural to consider constructing a probability path for two graph distributions with MRFs as the backbone.

A.2 Derivation of graph Markov random fields

We show the derivation of Definition 2, which is restated here:

Definition 3 (Graph Markov Random Fields). GraphMRF statistically describes graphs as,

$$p(\mathcal{G}; \mathbf{G}) = p(\mathcal{X}, \mathcal{E}; \mathbf{X}, \mathbf{W}) = p(\mathcal{X}; \mathbf{X}, \mathbf{W}) \cdot p(\mathcal{E}; \mathbf{W}) \text{ where } \mathcal{E} \sim \delta(\mathbf{W}) \text{ and} \quad (16)$$

$$\text{vec}(\mathcal{X}) \sim \mathcal{N}(\mathbf{X}, \mathbf{\Lambda}^\dagger), \text{ with } \mathbf{X} = \text{vec}(\mathbf{V}^\dagger \boldsymbol{\mu}), \mathbf{\Lambda} = (\nu \mathbf{I} + \mathbf{L}) \otimes \mathbf{V}^\top \mathbf{V}.$$

The \otimes is the Kronecker product, $\text{vec}(\cdot)$ is the vectorization operator and \mathbf{I} is the identity matrix.

Derivation:

We start from

$$p(\mathcal{X}; \mathbf{X}, \mathbf{W}) \propto \prod_v \exp\{-(\nu + d_v) \|\mathbf{V} x_v - \boldsymbol{\mu}_v\|^2\} \prod_{u,v} \exp\{w_{uv} [(\mathbf{V} x_u - \boldsymbol{\mu}_u)^\top (\mathbf{V} x_v - \boldsymbol{\mu}_v)]\}. \quad (17)$$

We assume that the linear transformation matrix has dimension $\mathbf{V} \in \mathbb{R}^{K' \times K}$ given that $x_v \in \mathbb{R}^K$ and define a transformed variable

$$h_v \equiv \mathbf{V} x_v - \boldsymbol{\mu}_v \in \mathbb{R}^{K'}, \text{ stacking as } \mathcal{H} \in \mathbb{R}^{|\mathcal{V}| \times K'}. \quad (18)$$

The probability becomes

$$P(\mathcal{H}; \mathbf{X}, \mathbf{W}) \propto \prod_v \exp\{-(\nu + d_v) \|h_v\|^2\} \prod_{u,v} \exp\{w_{uv} h_u^\top h_v\}. \quad (19)$$

Then, the terms inside the exponent in Eq. (19) become

$$\begin{aligned} -\sum_v (\nu + d_v) \|h_v\|^2 + \sum_{u,v} w_{uv} h_u^\top h_v &= -\sum_v (\nu + d_v) h_v^\top h_v + \sum_{u,v} w_{uv} h_u^\top h_v \\ &= -\sum_{u,v} h_u^\top [(\nu + d_u) \delta_{uv} - w_{uv}] h_v, \end{aligned}$$

where the Kronecker delta $\delta_{uv} = 1$ if $u = v$ and 0 else. We define a squared matrix $\mathbf{\Lambda}'$ to arrange the inner term, which can be written as,

$$\mathbf{\Lambda}' = \nu \mathbf{I} + \mathbf{L} \quad \text{with} \quad \mathbf{\Lambda}'_{uv} = (\nu + d_u) \delta_{uv} - w_{uv}. \quad (20)$$

\mathbf{I} is the identity matrix. Thus, the exponent in compact matrix form gives

$$-\frac{1}{2} \text{Tr}(\mathcal{H}^\top \mathbf{\Lambda}' \mathcal{H}), \text{ where } \mathcal{H} = \begin{pmatrix} h_1 \\ h_2 \\ \vdots \\ h_{|\mathcal{V}|} \end{pmatrix}. \quad (21)$$

It is possible to rearrange the exponent as

$$\text{Tr}(\mathcal{H}^\top \mathbf{\Lambda}' \mathcal{H}) = \text{vec}(\mathcal{H})^\top (\mathbf{\Lambda} \otimes \mathbf{I}) \text{vec}(\mathcal{H}), \quad (22)$$

where \otimes denotes the Kronecker product. This is exactly in the form of a multivariate colored Gaussian. Thus, the joint distribution of $\text{vec}(\mathcal{H})$ (of dimension $|\mathcal{V}|K'$) is given by

$$\text{vec}(\mathcal{H}) \sim \mathcal{N}\left(0, (\nu \mathbf{I} + \mathbf{L})^{-1} \otimes \mathbf{I}_{K'}\right), \quad (23)$$

Recall that $h_v = \mathbf{V} x_v - \boldsymbol{\mu}_v$, we obtain

$$\text{vec}(\mathcal{H}) = (\mathbf{I} \otimes \mathbf{V}) \text{vec}(\mathcal{X}) - \text{vec}(\boldsymbol{\mu}). \quad (24)$$

Since the transformation is linear, the distribution over \mathcal{X} remains Gaussian. By the properties of linear transformations of Gaussians, if

$$\text{vec}(\mathcal{H}) \sim \mathcal{N}(\text{vec}(\boldsymbol{\mu}), \boldsymbol{\Sigma}_h), \text{vec}(\mathcal{X}) = (\mathbf{I} \otimes \mathbf{V}^\dagger) \text{vec}(\mathcal{H}), \quad (25)$$

then

$$\text{vec}(\mathcal{X}) \sim \mathcal{N}\left((\mathbf{I} \otimes \mathbf{V}^\dagger) \text{vec}(\boldsymbol{\mu}), (\mathbf{I}_n \otimes \mathbf{V}^\dagger) \boldsymbol{\Sigma}_h (\mathbf{I}_n \otimes \mathbf{V}^\dagger)^\top\right). \quad (26)$$

Thus, using the mixed-product property of the Kronecker product,

$$(\mathbf{I} \otimes \mathbf{V}^\dagger)((\nu \mathbf{I} + \mathbf{L})^{-1} \otimes \mathbf{I})(\mathbf{I}_n \otimes \mathbf{V}^\dagger)^\top = (\mathbf{L} + \nu \mathbf{I})^{-1} \otimes (\mathbf{V}^\dagger \mathbf{V}^{\dagger\top}) \quad (27)$$

Finally, the joint distribution over \mathcal{X} is

$$\begin{aligned} \text{vec}(\mathcal{X}) &\sim \mathcal{N}(\mathbf{X}, \boldsymbol{\Sigma}), \\ \text{with } \mathbf{X} &= (\mathbf{I} \otimes \mathbf{V}^\dagger) \text{vec}(\boldsymbol{\mu}) = \text{vec}(\mathbf{V}^\dagger \boldsymbol{\mu}) \\ \text{and } \boldsymbol{\Sigma} &= (\nu \mathbf{I} + \mathbf{L})^{-1} \otimes (\mathbf{V}^\dagger \mathbf{V}^{\dagger\top}), \end{aligned} \quad (28)$$

We use the following lemma:

Lemma 1. *Given two invertible matrices \mathbf{A} and \mathbf{B} , their Kronecker product satisfies $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$.*

So that we get

$$\text{vec}(\mathcal{X}) \sim \mathcal{N}(\mathbf{X}, \boldsymbol{\Lambda}^\dagger), \text{ with } \mathbf{X} = \text{vec}(\mathbf{V}^\dagger \boldsymbol{\mu}), \boldsymbol{\Lambda} = (\nu \mathbf{I} + \mathbf{L}) \otimes \mathbf{V}^\top \mathbf{V}. \quad (29)$$

which ends the derivation.

A.3 The usage scope of graph Markov random fields

Given that our Graph Markov random fields (GraphMRF) have an explicit form to constrain the graph distribution, it inherits certain inductive biases and we have to properly understand the usage scenarios.

To understand the scenarios which GraphMRF could be used, we start by stating the plain form of MRF by melting the linear transformer matrix and the mean term, i.e., giving $h_v = \mathbf{V} x_v - \boldsymbol{\mu}$ which gives the probability density as,

$$P(\mathcal{H} | \mathbf{L}) \propto \exp(-\text{trace}(\mathcal{H}^\top (\mathbf{L} + \nu \mathbf{I}) \mathcal{H})) = \exp(-\mathbf{w}_{uv} \|h_u - h_v\|^2 - \nu \sum_u \|h_u\|_F^2) \quad (30)$$

Plain graph Markov random fields. To understand the scenarios which we can utilize MRF to model graphs, we first consider the simplest case when V is rectangular orthogonal (semi-orthogonal) matrix such that $V^T V = I$ and the mean $\mu = 0$, the probability density becomes,

$$P(X, L) \propto \exp(-X^T(L + \nu I)X) = \exp(-\sum_{\{u,v\} \in \mathcal{E}} W_{uv}(\mathbf{x}_u - \mathbf{x}_v)^2 - \nu \sum_u \mathbf{x}_u^2) \quad (31)$$

As $\nu \rightarrow 0$, the exponent term inside becomes

$$S(X, L) = -\sum_{\{u,v\} \in \mathcal{E}} W_{uv}(\mathbf{x}_u - \mathbf{x}_v)^2 = \text{trace}(X^T L X), \quad (32)$$

where we name $S(X, L)$ as the smoothness of the graph features. The smoothness measures how similar the neighbors connected are. For instance, if there exists an edge between node u and v weighted as W_{uv} , the likelihood will be higher if \mathbf{x}_u and \mathbf{x}_v be similar, so that $\|\mathbf{x}_u - \mathbf{x}_v\|^2$ are small. This suggests that the probability will be higher if the $S(X, L)$ is small.

This vanilla form captures the first type of graphs the GraphMRF can model - the homophily graphs, i.e., similar nodes (measured by the node attribute) may be more likely to attach to each other than dissimilar ones. This includes the social networks with friendship, traffic networks, etc.

Graph Markov random fields with embeddings. Now we move one step further to consider the graphs with linear transformer matrix V . Linear transformation provides a map from the feature space to the latent space, which can be considered as an embedding method to empower the models with better expressiveness. As a simple example, when the V provides a negative projection, the mapping can capture the heterophily relationships, which means the nodes connected are dissimilar.

Coincidentally, this aligns well with the famous embedding method Node2Vec as in Grover and Leskovec [21], where the edge weights are proportional to the negative distance, or the inner product of the embeddings. i.e.,

$$W_{uv} \propto \exp(-\|V\mathbf{x}_u - V\mathbf{x}_v\|_F^2) \quad (33)$$

In [27] it is derived that learning the parameters of MRFs is intrinsically equivalent to learning embeddings similar to Node2Vec. As such, the expressiveness of MRFs are as good as Node2Vec, which grants its usage to molecule graphs, protein interaction networks, social networks, and knowledge graphs. In our paper we make the assumption is that the linear mapping from X the observation is shared. This requirement translates to that the two graphs should have the same embedding space and feature space, which is practical if the reference distribution and data distributions share the same space.

Graphs without features. We wish to emphasize that even though the GraphMRF is constructed under the assumption that graph features exist, it is capable of modeling the non-attributed graphs, such as planar and SBM graphs. To do so, we consider the optimization over the Rayleigh function: It is shown that, if v_1, \dots, v_{k-1} are orthonormal eigenvectors for $\lambda_1, \dots, \lambda_{k-1}$, then the eigenvalues satisfy,

$$\lambda_k = \min_{\substack{\mathbf{x} \neq 0 \\ \mathbf{x} \perp v_1, \dots, v_{k-1}}} R(\mathbf{x}), \text{ with } R(\mathbf{x}) = \frac{\mathbf{x}^T L \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \quad (34)$$

In such a scenario, the graphs are no longer related to the actual node features, but instead, the eigenvectors v_k serve as an intrinsic graph feature. Interestingly, $R(\mathbf{x})$ is the normalized form of the smoothness as in Eq. (31), which is exactly equivalent to the eigenvalues. This means that if the spectrum of the graph Laplacian focuses on the low-frequency components, the corresponding graphs would give a higher probability in terms of MRFs if we view the eigenvectors as the node features. Such a pattern could capture the graph distributions almost all the synthetic datasets, such as Planar, SBM, TLS, COMM20 datasets, satisfy such a property. There is no surprise that our models capture better global patterns in those datasets. But it is also worth pointing out that there exists exceptions, such as the tree graph, which does not have a clear clustering pattern, thus the spectrum does not follow our GraphMRF.

Future work. A limitation of our method is that it cannot easily capture the generation of graphs with multiple relation types, which we name heterogeneous graphs. Even though we utilize an intuitive solution in the experiment to produce Table 2: we first sample the pure graph structure without

edge types to produce the graph backbone, and then sample the edge types via liner interpolated probability on top of the backbone. The solution provides preliminary results for the graph generation in multi-relational graphs, but still requires improvements. Fortunately, there exists a few ways to extend the GraphMRF to heterogeneous graphs [27]. An interesting future work can be generalizing our model to heterogeneous graphs by considering GraphMRF variants, such as the H2MN proposed in Jiang et al. [27].

B Proofs

B.1 Wasserstein Distance between two colored gaussian distributions

We first prove the lemma that captures the Wasserstein distance between two colored Gaussians, which will be used in deriving our Bures-Wasserstein distances in graph generations.

Lemma 2. *Consider two measures $\eta_0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0)$ and $\eta_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1)$, describing two colored Gaussian distributions with mean $\boldsymbol{\mu}_0, \boldsymbol{\mu}_1$ and covariance matrices Σ_0, Σ_1 . Then the Wasserstein distance between these probability distributions is given by*

$$(\mathcal{W}_2(\eta_0, \eta_1))^2 = \|\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1\|^2 + \text{Tr}\left(\Sigma_0 + \Sigma_1 - 2\left(\Sigma_0^{1/2}\Sigma_1\Sigma_0^{1/2}\right)^{1/2}\right).$$

Proof. We first state the following proposition.

Proposition 4. *(Translation Invariance of the 2-Wasserstein Distance for Gaussian Measures) Consider two measures $\eta_0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0)$ and $\eta_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1)$ and their centered measure as $\tilde{\eta}_0 = \mathcal{N}(0, \Sigma_0)$ and $\tilde{\eta}_1 = \mathcal{N}(0, \Sigma_1)$, the squared Wasserstein distance decomposes as*

$$\mathcal{W}_2^2(\eta_0, \eta_1) = \|\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1\|_2^2 + \mathcal{W}_2^2(\tilde{\eta}_0, \tilde{\eta}_1)$$

Proof:

Consider two random vectors \mathcal{X}, \mathcal{Y} distributed as η_0, η_1 ,

$$\mathcal{X} = \boldsymbol{\mu}_0 + \tilde{\mathcal{X}}, \mathcal{Y} = \boldsymbol{\mu}_1 + \tilde{\mathcal{Y}}, \text{ with } \tilde{\mathcal{X}} \sim \tilde{\eta}_0, \tilde{\mathcal{Y}} \sim \tilde{\eta}_1.$$

For any coupling $(\mathcal{X}, \mathcal{Y})$, we consider the expected squared Euclidean distance,

$$\begin{aligned} \mathbb{E}_{\mathcal{X}, \mathcal{Y}} \|\mathcal{X} - \mathcal{Y}\|^2 &= \mathbb{E}_{\mathcal{X}, \mathcal{Y}} \|\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1 + (\tilde{\mathcal{X}} - \tilde{\mathcal{Y}})\|^2 \\ &= \|\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1\|^2 + 2\langle \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1, \tilde{\mathcal{X}} - \tilde{\mathcal{Y}} \rangle + \mathbb{E}_{\tilde{\mathcal{X}}, \tilde{\mathcal{Y}}} \|\tilde{\mathcal{X}} - \tilde{\mathcal{Y}}\|^2 \end{aligned} \quad (35)$$

Since $\tilde{\mathcal{X}}$ and $\tilde{\mathcal{Y}}$ both have zero mean, we have $\mathbb{E}[\tilde{\mathcal{X}} - \tilde{\mathcal{Y}}] = 0$ so the cross-term vanishes. Thus,

$$\mathbb{E} \|\mathcal{X} - \mathcal{Y}\|^2 = \|\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1\|^2 + \mathbb{E} \|\tilde{\mathcal{X}} - \tilde{\mathcal{Y}}\|^2 \quad (36)$$

Take the definition of 2-Wasserstein distance, the infimum over all couplings directly yields

$$\begin{aligned} (\mathcal{W}_2(\eta_0, \eta_1))^2 &= \inf_{\pi \in \Pi(\eta_0, \eta_1)} \int \|\mathcal{X} - \mathcal{Y}\|^2 d\pi(\mathcal{X}, \mathcal{Y}). \\ &= \|\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1\|^2 + \mathcal{W}_2^2(\tilde{\eta}_0, \tilde{\eta}_1) \end{aligned} \quad (37)$$

This completes the proof of Proposition 4.

Now we prove the flowing proposition, which will give us our lemma.

Proposition 5. *Given two centered measures as $\tilde{\eta}_0 = \mathcal{N}(0, \Sigma_0)$ and $\tilde{\eta}_1 = \mathcal{N}(0, \Sigma_1)$*

$$\mathcal{W}_2^2(\tilde{\eta}_0, \tilde{\eta}_1) = \text{Tr}\left(\Sigma_0 + \Sigma_1 - 2\left(\Sigma_0^{1/2}\Sigma_1\Sigma_0^{1/2}\right)^{1/2}\right). \quad (38)$$

proof. The coupling π of $\tilde{\eta}_0$ and $\tilde{\eta}_1$ is a joint Gaussian measure with zero mean and covariance matrix

$$\Sigma_c = \begin{pmatrix} \Sigma_0 & C \\ C^T & \Sigma_1 \end{pmatrix} \geq 0, \quad (39)$$

where C is the cross-covariance and \geq means the matrix is positive semi-definitive (PSD). The expected squared distance between the two random vectors $(\mathcal{X}, \mathcal{Y})$ drawn from π is then described as,

$$\begin{aligned}\mathbb{E}\|\mathcal{X} - \mathcal{Y}\|^2 &= \text{Tr}(\mathbb{E}[(\mathcal{X} - \mathcal{Y})(\mathcal{X} - \mathcal{Y})^\top]) \\ &= \text{Tr}(\Sigma_0) + \text{Tr}(\Sigma_1) - 2 \text{Tr}(C).\end{aligned}\quad (40)$$

The definition of Wasserstein distance gives,

$$\mathcal{W}_c(\eta_0, \eta_1) = \inf_{\pi \in \Pi(\eta_0, \eta_1)} \mathbb{E}\|\mathcal{X} - \mathcal{Y}\|^2 \quad (41)$$

Thus, minimizing the Wasserstein distance is equivalent to maximizing $\text{Tr}(C)$ over all C subject to the joint covariance is positive semi-definite (PSD). It turns out (see Dowson and Landau [16], Olkin and Pukelsheim [42], Takatsu [52]) that the condition in Eq. (39) is equivalent to,

$$\Sigma_1 - C^\top \Sigma_0^{-1} C \geq 0 \leftrightarrow \Sigma_0^{-1/2} C \Sigma_1^{-1/2} \text{ has operator norm } \leq 1 \quad (42)$$

So we denote $K := \Sigma_0^{-1/2} C \Sigma_1^{-1/2}$ with $\|K\|_{\text{op}} \leq 1$. Then

$$\text{Tr}(C) = \text{Tr}(\Sigma_0^{1/2} K \Sigma_1^{1/2}) = \text{Tr}(K \Sigma_1^{1/2} \Sigma_0^{1/2}).$$

Using von Neumann trace inequality, its trace inner-product with K is maximized by choosing $K = I$ on the support.

$$\max_{\|K\|_{\text{op}} \leq 1} \text{Tr}(KA) = \text{Tr}(M^{1/2}), \quad M = \sqrt{AA^\top} = \Sigma_1^{1/2} \Sigma_0 \Sigma_1^{1/2}$$

Hence the optimal value of $\text{Tr}(C)$ is

$$\text{Tr}(C^*) = \text{Tr}\left[\left(\Sigma_1^{1/2} \Sigma_0 \Sigma_1^{1/2}\right)^{1/2}\right]$$

Substituting this optimal value into the expression of Wasserstein distance, we obtain

$$\mathcal{W}_2^2(\tilde{\eta}_0, \tilde{\eta}_1) = \text{Tr}(\Sigma_0) + \text{Tr}(\Sigma_1) - 2 \text{Tr}\left[\left(\Sigma_1^{1/2} \Sigma_0 \Sigma_1^{1/2}\right)^{1/2}\right]. \quad (43)$$

This completes the proof of proposition 5. Taking Proposition 4 and Proposition 5 together, we proved Lemma 2.

B.2 Derivation of the graph Wasserstein distance under MRF

We then prove the Bures-Wasserstein distance for two graph distributions. We restate Proposition 1,

Proposition 6 (Bures-Wasserstein Distance). *Consider two same-sized graphs $\mathcal{G}_0 \sim p(\mathcal{X}_0, \mathcal{E}_0)$ and $\mathcal{G}_1 \sim p(\mathcal{X}_1, \mathcal{E}_1)$ with V shared for two graphs, described by the distribution in Definition 2. When the graphs are equipped with graph Laplacian matrices L_0 and L_1 satisfying 1) is Positive Semi-Definite (PSD) and 2) has only one zero eigenvalue. The Bures-Wasserstein distance between these two random graph distributions is given by*

$$d_{BW}(\mathcal{G}_0, \mathcal{G}_1) = \|\mathbf{X}_0 - \mathbf{X}_1\|_F^2 + \beta \text{Tr}\left(L_0^\dagger + L_1^\dagger - 2\left(L_0^{\dagger/2} L_1^\dagger L_0^{\dagger/2}\right)^{1/2}\right), \quad (44)$$

as $\nu \rightarrow 0$ and β is a constant related to the norm of V .

Specifically, Definition 2 uses graph Markov random fields to describe a graph as

$$\begin{aligned}p(\mathcal{G}; \mathbf{G}) &= p(\mathcal{X}, \mathcal{E}; \mathbf{X}, \mathbf{W}) = p(\mathcal{X}; \mathbf{X}, \mathbf{W}) \cdot p(\mathcal{E}; \mathbf{W}) \text{ where } \mathcal{E} \sim \delta(\mathbf{W}) \text{ and} \\ \text{vec}(\mathcal{X}) &\sim \mathcal{N}(\mathbf{X}, \mathbf{A}^\dagger), \text{ with } \mathbf{X} = \text{vec}(V^\dagger \boldsymbol{\mu}), \mathbf{A} = (\nu \mathbf{I} + L) \otimes V^\top V.\end{aligned}\quad (45)$$

With the graph Wasserstein distance defined as,

$$(\text{Graph Wasserstein Distance}) \quad d_{BW}(\mathcal{G}_0, \mathcal{G}_1) := \mathcal{W}_c(\eta_{\mathcal{G}_0}, \eta_{\mathcal{G}_1}) = \mathcal{W}_c(\eta_{\mathcal{X}_0}, \eta_{\mathcal{X}_1}) + \mathcal{W}_c(\eta_{\mathcal{E}_0}, \eta_{\mathcal{E}_1}).$$

We first consider calculating $\mathcal{W}_c(\eta_{\mathcal{X}_0}, \eta_{\mathcal{X}_1})$. Specifically, this is the distance between two colored Gaussian measures where

$$\eta_i \sim \mathcal{N}(\boldsymbol{\mu}'_i, \Sigma_i), \quad i = 0, 1, \quad (46)$$

where $\boldsymbol{\mu}'_i = V_i \otimes \boldsymbol{\mu}_i$ and $\Sigma_i^{-1} = \Lambda_i = (\nu \mathbf{I} + L_i) \otimes (V_i^\top V_i)$.

where we first assume that these two Gaussians are emitted from different linear transformation matrices \mathbf{V}_0 and \mathbf{V}_1 . This will bring us the most general and flexible form that could be universally applicable, and potentially can bring more insights to future work. Next, we will inject a few assumptions to arrive at a more practical form for building the flow matching models.

An important property of Kronecker product: Given two invertible matrices \mathbf{A} and \mathbf{B} , their Kronecker product satisfies $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$. Using such a property, in the limit as $\nu \rightarrow 0$, we have

$$\Lambda_i \rightarrow \mathbf{L}_i \otimes (\mathbf{V}_i^\top \mathbf{V}_i) \implies \Sigma_i = \mathbf{L}_i^{-1} \otimes (\mathbf{V}_i^\top \mathbf{V}_i)^{-1}. \quad (47)$$

According to Propensity 2, the squared 2-Wasserstein distance between two Gaussian measures is given by

$$\mathcal{W}_2^2(\eta_0, \eta_1) = \underbrace{\|\boldsymbol{\mu}'_0 - \boldsymbol{\mu}'_1\|^2}_{\text{Mean term}} + \underbrace{\text{Tr}\left(\Sigma_0 + \Sigma_1 - 2\left(\Sigma_0^{1/2}\Sigma_1\Sigma_0^{1/2}\right)^{1/2}\right)}_{\text{Covariance Term}}. \quad (48)$$

Mean Term. Since $\boldsymbol{\mu}'_i = \mathbf{V} \otimes \boldsymbol{\mu}_i$, the mean difference becomes

$$\|\boldsymbol{\mu}'_0 - \boldsymbol{\mu}'_1\|^2 = \|\mathbf{V}_0\boldsymbol{\mu}_0 - \mathbf{V}_1\boldsymbol{\mu}_1\|_F^2 = \|\mathbf{X}_0 - \mathbf{X}_1\|_F^2 \quad (49)$$

Covariance term. Using the property of the Kronecker product, the square root of Eq. (47) factors in as

$$\Sigma_i^{1/2} = \mathbf{L}_i^{-1/2} \otimes (\mathbf{V}_i^\top \mathbf{V}_i)^{-1/2}. \quad (50)$$

and

$$\Sigma_0^{1/2}\Sigma_1\Sigma_0^{1/2} = \left(\mathbf{L}_0^{-1/2}\mathbf{L}_1^{-1}\mathbf{L}_0^{-1/2}\right) \otimes \left((\mathbf{V}_0^\top \mathbf{V}_0)^{-1/2}(\mathbf{V}_1^\top \mathbf{V}_1)^{-1}(\mathbf{V}_0^\top \mathbf{V}_0)^{-1/2}\right) \quad (51)$$

We first look into the term related to \mathbf{V}_0 and \mathbf{V}_1 , which is,

$$\begin{aligned} \text{Tr}\left((\mathbf{V}_0^\top \mathbf{V}_0)^{-1/2}(\mathbf{V}_1^\top \mathbf{V}_1)^{-1}(\mathbf{V}_0^\top \mathbf{V}_0)^{-1/2}\right) &= \text{Tr}\left((\mathbf{V}_1^\top \mathbf{V}_1)^{-1}(\mathbf{V}_0^\top \mathbf{V}_0)^{-1/2}(\mathbf{V}_0^\top \mathbf{V}_0)^{-1/2}\right) \\ &= \text{Tr}\left((\mathbf{V}_1^\top \mathbf{V}_1)^{-1}(\mathbf{V}_0^\top \mathbf{V}_0)^{-1}\right) \end{aligned} \quad (52)$$

As $\text{Tr}(\mathbf{A} + \mathbf{B}) = \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B})$ the covariance term becomes

Covariance Term

$$\begin{aligned} &= \text{Tr}\left(\Sigma_0 + \Sigma_1 - 2\left(\Sigma_0^{1/2}\Sigma_1\Sigma_0^{1/2}\right)^{1/2}\right) \\ &= \text{Tr}(\Sigma_0) + \text{Tr}(\Sigma_1) - 2\text{Tr}\left((\Sigma_0^{1/2}\Sigma_1\Sigma_0^{1/2})^{1/2}\right) \\ &= \text{Tr}\left(\mathbf{L}_0^{-1} \otimes (\mathbf{V}_0^\top \mathbf{V}_0)^{-1} + \mathbf{L}_1^{-1} \otimes (\mathbf{V}_1^\top \mathbf{V}_1)^{-1} - 2\left(\mathbf{L}_0^{-1/2}\mathbf{L}_1^{-1}\mathbf{L}_0^{-1/2}\right)^{1/2} \otimes (\mathbf{V}_1^\top \mathbf{V}_1)^{-1/2}(\mathbf{V}_0^\top \mathbf{V}_0)^{-1/2}\right) \end{aligned} \quad (53)$$

Given that $\text{Tr}(\mathbf{A} \otimes \mathbf{B}) = \text{Tr}(\mathbf{A})\text{Tr}(\mathbf{B})$ and $\text{Tr}(\mathbf{V}^\top \mathbf{V}) = \|\mathbf{V}\|_F^2$ for any real-valued matrix \mathbf{V} , we can further derive,

$$\begin{aligned} \text{Covariance Term} &= \text{Tr}[(\mathbf{V}_0^\top \mathbf{V}_0)^{-1}]\text{Tr}(\mathbf{L}_0^\dagger) + \text{Tr}[(\mathbf{V}_1^\top \mathbf{V}_1)^{-1}]\text{Tr}(\mathbf{L}_1^\dagger) \\ &\quad - 2\text{Tr}\left(\mathbf{L}_0^{\dagger/2}\mathbf{L}_1^\dagger\mathbf{L}_0^{\dagger/2}\right)^{1/2} \cdot \text{Tr}[(\mathbf{V}_1^\top \mathbf{V}_1)^{-1/2}(\mathbf{V}_0^\top \mathbf{V}_0)^{-1/2}]. \end{aligned} \quad (54)$$

Unfortunately, to simplify this equation, we have to make the two gram matrix, $(\mathbf{V}_0^\top \mathbf{V}_0)^{-1}$ and $(\mathbf{V}_1^\top \mathbf{V}_1)^{-1}$ agree, i.e., $(\mathbf{V}_1^\top \mathbf{V}_1)^{-1} = (\mathbf{V}_0^\top \mathbf{V}_0)^{-1}$. This will be satisfied if and only if there exists an orthogonal matrix \mathbf{Q} such that

$$\mathbf{V}_1^\dagger = \mathbf{V}_0^\dagger \mathbf{Q}.$$

Thus, to further process, we simply consider the case when \mathbf{V}_1 and \mathbf{V}_0 are exactly the same, i.e., $\mathbf{V}_1 = \mathbf{V}_0 = \mathbf{V}$ (we have already discussed how realistic this assumption is in Appendix A.3). So that we work under the assumptions that $\|\mathbf{V}_0^\dagger\|_F^2 = \|\mathbf{V}_1^\dagger\|_F^2 = \beta$, which simplify the trace as

$$\text{Covariance Term} = \beta \cdot \text{Tr}\left(\mathbf{L}_0^\dagger + \mathbf{L}_1^\dagger - 2\left(\mathbf{L}_0^{\dagger/2}\mathbf{L}_1^\dagger\mathbf{L}_0^{\dagger/2}\right)^{1/2}\right). \quad (55)$$

Combining the mean term and the covariance term, we obtain the Wasserstein distance of $\mathcal{W}_c(\eta_{\mathcal{X}_0}, \eta_{\mathcal{X}_1})$

For calculating $\mathcal{W}_c(\eta_{\mathcal{E}_0}, \eta_{\mathcal{E}_1})$, we have the freedom to choose the cost function when obtaining the Wasserstein distance. Note that \mathbf{W} serves as the prior for the Gaussian covariance matrix Σ , where the covariance has to be positive-semi definite. Thus, according to [5], a proper distance between two positive semi-definite matrices is measured by

$$\mathcal{W}(\eta_{\mathcal{E}_0}, \eta_{\mathcal{E}_1}) = \left\| \Sigma_0^{1/2} - \Sigma_1^{1/2} \right\|_F^2. \quad (56)$$

Coincidentally, this is another usage case when the Bures-Wasserstein metric is utilized. Putting everything together, the Wasserstein distance in the limit $\nu \rightarrow 0$ is

$$\begin{aligned} d_{\text{BW}}(\mathcal{G}_0, \mathcal{G}_1) &= \|\mathbf{V}_0 \boldsymbol{\mu}_0 - \mathbf{V}_1 \boldsymbol{\mu}_1\|_F^2 + (\beta + 1) \cdot \text{Tr} \left(\mathbf{L}_0^\dagger + \mathbf{L}_1^\dagger - 2 \left(\mathbf{L}_0^{\dagger/2} \mathbf{L}_1^\dagger \mathbf{L}_0^{\dagger/2} \right)^{1/2} \right). \\ &= \underbrace{\|\mathbf{X}_0 - \mathbf{X}_1\|_F^2}_{d_{\mathbf{X}}(\mathbf{X}_0, \mathbf{X}_1)} + (\beta + 1) \cdot \underbrace{\text{Tr} \left(\mathbf{L}_0^\dagger + \mathbf{L}_1^\dagger - 2 \left(\mathbf{L}_0^{\dagger/2} \mathbf{L}_1^\dagger \mathbf{L}_0^{\dagger/2} \right)^{1/2} \right)}_{d_{\mathbf{L}}(\mathbf{L}_0, \mathbf{L}_1)}. \end{aligned} \quad (57)$$

This expression separates the contribution of the mean difference (transformed by \mathbf{V}) and the discrepancy between the covariance structures (encoded in \mathbf{L}_0 and \mathbf{L}_1). This could be further used to derive BW interpolation, which we will show in Appendix C.1. In the main body, constant β actually corresponds to $\beta + 1$ here. This complete our derivation in Proposition 1.

C Derivation of Bures-Wasserstein flow matching

In order to build the flow matching framework, we need to derive the optimal interpolation and the corresponding velocities for the probability path $p(G_t | G_0, G_1)$. This is achieved via the OT displacement between two graph distributions.

C.1 The Bures-Wasserstein graph interpolation

We aim to recover the proposition stated as follows.

Proposition 7 (Bures-Wasserstein interpolation). *The graph minimizer of Eq. (10), $\mathcal{G}_t = \{\mathcal{V}, \mathcal{E}_t, \mathcal{X}_t\}$, have its node features following a colored Gaussian distribution, $\mathcal{X}_t \sim \mathcal{N}(\mathbf{X}_t, \boldsymbol{\Lambda}_t^\dagger)$ with $\boldsymbol{\Lambda}_t = (\nu \mathbf{I} + \mathbf{L}_t) \otimes \mathbf{V}^\top \mathbf{V}$ and edges following $\mathcal{E}_t \sim \delta(\mathbf{W}_t)$, specifically,*

$$\mathbf{L}_t^\dagger = \mathbf{L}_0^{1/2} \left((1-t) \mathbf{L}_0^\dagger + t \left(\mathbf{L}_0^{\dagger/2} \mathbf{L}_1^\dagger \mathbf{L}_0^{\dagger/2} \right)^{1/2} \right)^2 \mathbf{L}_0^{1/2}, \quad \mathbf{X}_t = (1-t) \mathbf{X}_0 + t \mathbf{X}_1 \quad (58)$$

The interpolation is an extension of the concept of mean, where in the optimal transport world, the Wasserstein barycenter (mean) of measures $\eta_0, \dots, \eta_{m-1}$ under weights $\lambda_0, \dots, \lambda_{m-1}$ can be derived over the following optimization problem:

$$\bar{\eta} = \arg \min_{\eta} \sum_{j=0}^{m-1} \lambda_j (\mathcal{W}_2(\eta, \eta_j))^2 \quad (59)$$

When $m = 2$, based on the Bures-Wasserstein (BW) distance, we can define the OT displacement minimization problem on graphs described as,

$$\mathcal{G}_t = \arg \min_{\tilde{\mathcal{G}}} (1-t) d_{\text{BW}}(\mathcal{G}_0, \tilde{\mathcal{G}}) + t d_{\text{BW}}(\tilde{\mathcal{G}}, \mathcal{G}_1). \quad (60)$$

where $d_{\text{BW}}(\mathcal{G}_0, \mathcal{G}_1)$ is described in Proposition 1. The optimal graph interpolation is the solution to the problem.

In the setting of graph, this becomes a two-variable optimization problem, where

$$\mathcal{X}_t, \mathcal{E}_t = \arg \min_{\tilde{\mathcal{X}}, \tilde{\mathcal{E}}} (1-t) d_{\text{BW}}(\mathcal{G}_0, \tilde{\mathcal{G}}) + t d_{\text{BW}}(\tilde{\mathcal{G}}, \mathcal{G}_1). \quad (61)$$

Fortunately, recall in Eq. (57) that our distance measurement $d_{\text{BW}}(\mathcal{G}_0, \mathcal{G}_1)$ is decomposed into $d_{\mathbf{X}}(\mathbf{X}_0, \mathbf{X}_1)$ and $d_{\mathbf{L}}(\mathbf{L}_0, \mathbf{L}_1)$, then the optimization over node and edges are disentangleable into solving the two sub optimization problem,

$$\begin{aligned} \text{Sub-question 1: } \quad \bar{\mathbf{X}}_t &= \arg \min_{\tilde{\mathbf{X}}} (1-t)\|\mathbf{X}_0 - \tilde{\mathbf{X}}\|_F^2 + t\|\tilde{\mathbf{X}} - \mathbf{X}_1\|_F^2 \\ \text{Sub-question 2: } \quad \bar{\mathbf{L}}_t &= \arg \min_{\tilde{\mathbf{L}}} (1-t)d_{\mathbf{L}}(\mathbf{L}_0, \tilde{\mathbf{L}}) + td_{\mathbf{L}}(\mathbf{L}_1, \tilde{\mathbf{L}}) \end{aligned} \quad (62)$$

This two problems are completely disentangled thus we can solve them separately.

Sub-question 1 For the first problem, we simply set the derivate to 0 and get,

$$(1-t)(\tilde{\mathbf{X}} - \mathbf{X}_0) + t(\tilde{\mathbf{X}} - \mathbf{X}_1) = 0 \rightarrow \mathbf{X}_t = (1-t)\mathbf{X}_0 + t(\mathbf{X}_1) \quad (63)$$

Subquestion 2 The second subproblem is equivalent in deriving the covariance of Bures-Wasserstein interpolation between two Gaussian measures, $\eta_0 \sim \mathcal{N}(0, \mathbf{L}_0^\dagger)$ and $\eta_1 \sim \mathcal{N}(0, \mathbf{L}_1^\dagger)$. This problem has been properly addressed in Haasler and Frossard [22] and here we just verbose their results. For more details we refer the reader to Haasler and Frossard [22] for a further discussion.

The optimal transport geodesic between $\eta_0 \sim \mathcal{N}(0, \mathbf{L}_0^\dagger)$ and $\eta_1 \sim \mathcal{N}(0, \mathbf{L}_1^\dagger)$ is defined by $\eta_t = ((1-t)I + t\mathbf{T})_{\#}\eta_0$, where the symbol “ $\#$ ” denotes the push-forward of a measure by a mapping, \mathbf{T} is a linear map that satisfies $\mathbf{T}\mathbf{L}_0^\dagger\mathbf{T}^\top = \mathbf{L}_1^\dagger$.

We define a new matrix \mathbf{M} and do normalization, which leads to,

$$\mathbf{T} = \mathbf{L}_0^{1/2} \mathbf{M} \mathbf{L}_0^{1/2} \quad (64)$$

Plug in gives,

$$\begin{aligned} \mathbf{T}\mathbf{L}_0^\dagger\mathbf{T}^\top &= \mathbf{L}_0^{1/2} \mathbf{M} \mathbf{L}_0^{1/2} \mathbf{L}_0^\dagger (\mathbf{L}_0^{1/2} \mathbf{M} \mathbf{L}_0^{1/2})^\top \\ &= \mathbf{L}_0^{1/2} \mathbf{M} \mathbf{M}^\top \mathbf{L}_0^{1/2}. \end{aligned} \quad (65)$$

So that we obtain

$$\mathbf{L}_1^\dagger = \mathbf{L}_0^{1/2} \mathbf{M} \mathbf{M}^\top \mathbf{L}_0^{1/2} \rightarrow \mathbf{M} = (\mathbf{L}_0^{\dagger/2} \mathbf{L}_1^\dagger \mathbf{L}_0^{\dagger/2})^{1/2} \quad (66)$$

Replace \mathbf{T} and we get,

$$\mathbf{T} = \mathbf{L}_0^{1/2} (\mathbf{L}_0^{\dagger/2} \mathbf{L}_1^\dagger \mathbf{L}_0^{\dagger/2})^{1/2} \mathbf{L}_0^{1/2} \quad (67)$$

Given that the geodesic $\eta_t = ((1-t)I + t\mathbf{T})_{\#}\eta_0$ which also has a Gaussian form $\eta_t \sim \mathcal{N}(0, \Sigma_t)$, We can then write the covariance matrix and obtain

$$\begin{aligned} \mathbf{L}_t^\dagger &= \Sigma_t = ((1-t)I + t\mathbf{T})_{\#}\mathbf{L}_0^\dagger((1-t)I + t\mathbf{T})^\top \\ &= \mathbf{L}_0^{1/2} \left((1-t)\mathbf{L}_0^\dagger + t \left(\mathbf{L}_0^{\dagger/2} \mathbf{L}_1^\dagger \mathbf{L}_0^{\dagger/2} \right)^{1/2} \right) \mathbf{L}_0^{1/2} \mathbf{L}_0^{1/2} \left((1-t)\mathbf{L}_0^\dagger + t \left(\mathbf{L}_0^{\dagger/2} \mathbf{L}_1^\dagger \mathbf{L}_0^{\dagger/2} \right)^{1/2} \right) \mathbf{L}_0^{1/2} \\ &= \mathbf{L}_0^{1/2} \left((1-t)\mathbf{L}_0^\dagger + t \left(\mathbf{L}_0^{\dagger/2} \mathbf{L}_1^\dagger \mathbf{L}_0^{\dagger/2} \right)^{1/2} \right)^2 \mathbf{L}_0^{1/2} \end{aligned} \quad (68)$$

Which ends the derivation.

Remark 1: Even though the GraphMRF in Definition 2 does rely on an implicit linear emission matrices \mathbf{V} , the BW interpolation in Proposition 2 can be obtained without explicitly accessing to the \mathbf{V} matrices. The property was attractive as in practice we can construct the probability path without explicitly fitting a \mathbf{V} beforehand.

C.2 Deriving the velocity of BW interpolation

We first show the general form of the velocity term for the Gaussian and Dirac measures.

Gaussian measure. For a time-parametrized Gaussian density $p_t(x) = \mathcal{N}(x; \boldsymbol{\mu}_t, \Sigma_t)$, the velocity field $v_t(x)$ satisfies the continuity equation

$$\partial_t p_t + \nabla \cdot (p_t v_t) = 0,$$

is an affine function of x . And the instantaneous velocity field follows,

$$v_t(\mathcal{X}) = \dot{\boldsymbol{\mu}}_t + \frac{1}{2} \dot{\Sigma}_t \Sigma_t^{-1} (\mathcal{X} - \boldsymbol{\mu}_t).$$

Dirac measure. When the measure is a Dirac function,

$$p_t(x) = \delta(\cdot, \mu_t).$$

We can just consider it as the limited case of the Gaussian measure, when $\Sigma_t \rightarrow 0$. So that the velocity at simply takes

$$v_t(x) = \dot{\mu}_t.$$

We then move to prove the following propensity for Bures-wasserstein velocity.

Proposition 8 (Bures-Wasserstein velocity). *For the graph \mathcal{G}_t following BW interpolation in Proposition 2, the conditional velocity at time t with observation G_t is given as,*

$$v_t(E_t | G_0, G_1) = \dot{\mathbf{W}}_t = \text{diag}(\dot{\mathbf{L}}_t) - \dot{\mathbf{L}}_t, \quad v_t(X_t | G_0, G_1) = \frac{1}{1-t}(\mathbf{X}_1 - \mathbf{X}_t) \quad (69)$$

$$\text{with } \dot{\mathbf{L}}_t = 2\mathbf{L}_t - \mathbf{T}\mathbf{L}_t - \mathbf{L}_t\mathbf{T} \text{ and } \mathbf{T} = \mathbf{L}_0^{1/2}(\mathbf{L}_0^\dagger \mathbf{L}_1^\dagger \mathbf{L}_0^\dagger)^{1/2} \mathbf{L}_0^{1/2}$$

where $\mathbf{W}_t = \mathbf{D}_t - \mathbf{L}_t$ and \mathbf{L}_t defined in Eq. (11).

Proof:

The graph structure velocity. As we assume the edges, $E_t \sim \delta(\cdot, \mathbf{W}_t)$, following a dirac distribution, the velocity is defined as

$$v_t(E_t) = \dot{\mathbf{W}}_t.$$

Given that, $\dot{\mathbf{W}}_t = \text{diag}(\dot{\mathbf{L}}_t) - \dot{\mathbf{L}}_t$, we transit to deriving the derivative of the Laplacian matrix, $\dot{\mathbf{L}}_t$. Using the fact that,

$$\frac{d}{dt}(\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \frac{d\mathbf{A}}{dt} \mathbf{A}^{-1}$$

we obtain the derivate of Laplacian matrix,

$$\dot{\mathbf{L}}_t = \frac{d(\Sigma_t^\dagger)}{dt} = \Sigma_t^\dagger \frac{d\Sigma_t}{dt} \Sigma_t^\dagger = \mathbf{L}_t \frac{d\Sigma_t}{dt} \mathbf{L}_t \quad (70)$$

According to Eq. (68) and Eq. (67), the covariance matrix is defined through the interpolation,

$$\Sigma_t = ((1-t)\mathbf{I} + t\mathbf{T})\mathbf{L}_0^\dagger((1-t)\mathbf{I} + t\mathbf{T}) := \mathbf{R}_t \mathbf{L}_0^\dagger \mathbf{R}_t \quad (71)$$

where $\mathbf{R}_t = (1-t)\mathbf{I} + t\mathbf{T}$. Taking the derivative, we get,

$$\dot{\Sigma}_t = \frac{d}{dt}(\mathbf{R}_t \Sigma_0 \mathbf{R}_t) = \mathbf{R}_t' \Sigma_0 \mathbf{R}_t + \mathbf{R}_t \Sigma_0 \mathbf{R}_t' = (\mathbf{T} - \mathbf{I}) \Sigma_0 \mathbf{R}_t + \mathbf{R}_t \Sigma_0 (\mathbf{T} - \mathbf{I}) \quad (72)$$

Using the fact that $\Sigma_0 \mathbf{R}_t = \mathbf{R}_t \Sigma_0 = \Sigma_t$, we obtain the covariance gradient

$$\dot{\Sigma}_t = (\mathbf{T} - \mathbf{I}) \Sigma_t + \Sigma_t (\mathbf{T} - \mathbf{I}) \quad (73)$$

So that,

$$\begin{aligned} -\dot{\mathbf{L}}_t &= \frac{d(\Sigma_t^\dagger)}{dt} = \Sigma_t^\dagger \frac{d\Sigma_t}{dt} \Sigma_t^\dagger = \mathbf{L}_t \frac{d\Sigma_t}{dt} \mathbf{L}_t \\ &= \mathbf{L}_t ((\mathbf{T} - \mathbf{I}) \mathbf{L}_t^\dagger + \mathbf{L}_t^\dagger (\mathbf{T} - \mathbf{I})) \mathbf{L}_t \\ &= \mathbf{L}_t (\mathbf{T} - \mathbf{I}) + (\mathbf{T} - \mathbf{I}) \mathbf{L}_t \\ &= \mathbf{L}_t \mathbf{T} + \mathbf{T} \mathbf{L}_t - 2\mathbf{L}_t \end{aligned} \quad (74)$$

Thus, $\dot{\mathbf{L}}_t = 2\mathbf{L}_t - \mathbf{L}_t \mathbf{T} - \mathbf{T} \mathbf{L}_t$.

Given that $\mathbf{L}_t = \mathbf{D}_t - \mathbf{W}_t$ so that $\mathbf{W}_t = \text{diag}(\mathbf{L}_t) - \mathbf{L}_t$, taking the derivative gives $\dot{\mathbf{W}}_t = \text{diag}(\dot{\mathbf{L}}_t) - \dot{\mathbf{L}}_t$. As we assume the edges, $E_t \sim \delta(\cdot, \mathbf{W}_t)$, the derivate directly yields the velocity,

$$v_t(E_t | G_0, G_1) = \dot{\mathbf{W}}_t = \text{diag}(\dot{\mathbf{L}}_t) - \dot{\mathbf{L}}_t.$$

The node feature velocity. The instantaneous velocity field follows,

$$v_t(\mathcal{X} \mid G_0, G_1) = \dot{\boldsymbol{\mu}}_t + \frac{1}{2} \dot{\Sigma}_t \Sigma_t^{-1} (\mathcal{X} - \boldsymbol{\mu}_t).$$

The mean gradient interpolating η_0 and η_1 can be written as $\dot{\boldsymbol{\mu}}_t = \mathbf{X}_1 - \mathbf{X}_0$ and $\mathbf{X}_t = (1-t)\mathbf{X}_0 + t\mathbf{X}_1$. So that the velocity leads to,

$$v_t(\mathcal{X} \mid G_0, G_1) = \mathbf{X}_1 - \mathbf{X}_0 + \frac{1}{2} \dot{\mathbf{L}}_t^\dagger \mathbf{L}_t (\mathcal{X} - \mathbf{X}_t).$$

However, in practice, we do not need such a complicated velocity term. We wish to avoid the estimation of complex gradient-inverse term so that we can escape from the complicated computation. Under the assumption that the amplitude of covariance is much smaller than the mean difference, we can omit the second term and just keep the mean difference. Hence the instantaneous velocity is simply described as

$$v_t(\mathbf{X}_t \mid G_0, G_1) = \mathbf{X}_1 - \mathbf{X}_0 = \mathbf{X}_1 - \mathbf{X}_0 = \frac{1}{1-t}(\mathbf{X}_1 - \mathbf{X}_t) \quad (75)$$

D Discrete Bures-Wasserstein flow matching for graph generation

D.1 Probability path construction for discrete Bures-Wasserstein flow matching

The discrete probability path. We design the probability path as discrete distributions,

$$\begin{aligned} p_t(x_v \mid G_0, G_1) &= \text{Categorical}([\mathbf{X}_t]_v), \quad p_t(e_{uv} \mid G_0, G_1) = \text{Bernoulli}([\mathbf{W}_t]_{uv}) \\ \text{s.t. } p_0(\mathcal{G}) &= \delta(G_0, \cdot), p_1(\mathcal{G}) = \delta(G_1, \cdot) \end{aligned} \quad (76)$$

where $\mathbf{W}_t = \mathbf{D}_t - \mathbf{L}_t$ with \mathbf{X}_t and \mathbf{L}_t defined the same in Eq. (11). We consider the fact that the Dirac distribution is a special case when the Categorical/Bernoulli distribution has probability 1 or 0, so the boundary condition $p_0(\mathcal{G}) = \delta(G_0, \cdot), p_1(\mathcal{G}) = \delta(G_1, \cdot)$ holds. As such, $\mathbf{X}_t = (1-t)\mathbf{X}_0 + t\mathbf{X}_1 \in [0, 1]^{|V| \times K}$. Since the boundary condition for each entry, $[\mathbf{X}_0]_v$ and $[\mathbf{X}_1]_v$ are two one-hot embeddings, $[\mathbf{X}_t]_v = t[\mathbf{X}_0]_v + (1-t)[\mathbf{X}_1]_v$ would sum to one, which works as a valid probability vector. Thus, $\text{Categorical}([\mathbf{X}_t]_v)$ is a K-class categorical distribution.

For the edge distribution, we just consider e_{uv} is conditionally independent of the other given $[\mathbf{W}_t]_{uv}$. One thing to emphasize is that, given the nature of Bures-Wasserstein interpolation, the yielded \mathbf{W}_t is not always bounded by $[0, 1]$ thus we have to hard-clip the boundary.

D.2 Approximating Wasserstein distance in Bernoulli distributions

To make sure that the individual nodes are structured and developed jointly while doing flow matching, we assume that the $\text{vec}(\mathcal{X})$ still maintains a covariance matrix similar to Eq. (8), which gives $\boldsymbol{\Lambda} = (\nu \mathbf{I} + \mathbf{L}) \otimes \mathbf{V}^\top \mathbf{V}$ given that \mathcal{X} is emitted from a latent variable \mathcal{H} through an affine transformation and the latent variable has a covariance matrix $(\nu \mathbf{I} + \mathbf{L})^{-1}$. Different from the Gaussian case, the latent variable would still be a discrete distribution, so that the affine transformation carries the covariance matrix out.

Unfortunately, the Wasserstein distance between two discrete graph distributions that follow Eq. (13) does not have a closed-form solution given the complex intertwined nature. However, it is possible to use the central limit theorem applied to \mathcal{X} so that we can approximate the Wasserstein distance of two Bernoulli distributions with the Gaussian counterpart. This approximation works when we are in high-dimensional case (high dimension means $|V|d$ is moderately large.), and the OT-distance between two such Bernoulli distributions is well-captured by the corresponding Gaussian formula, which we already introduced in Eq. (57).

With such nature, even though we are not sampling from Gaussian distributions anymore, it is possible to approximate the Wasserstein distance between two multivariate discrete distributions with the Gaussian counterpart, so the conclusions, such as optimal transport displacements, still hold. And we can similarly derive the Bures-Wasserstein velocity as in the next section.

D.3 Velocity for discrete Bures-Wasserstein flow matching

Node velocity. For node-wise, the path of node features \mathcal{X}_t can be re-written as $p_t(\mathcal{X}) = (1 - t)\delta(\cdot, \mathbf{X}_0) + t\delta(\cdot, \mathbf{X}_1)$ so the conditional velocity can be accessed through $v_t(X_t | G_0, G_1) = [\delta(\cdot, \mathbf{X}_1) - \delta(\cdot, \mathbf{X}_0)]/(1 - t)$ similar as the derivation in [19].

Edge velocity. For edge-wise, we look into each entry of the adjacency matrix \mathbf{W} , and consider a time-dependent Bernoulli distribution, the probability density function is:

$$p_t(e_{uv}) = [\mathbf{W}_t]_{uv}^{e_{uv}} (1 - [\mathbf{W}_t]_{uv})^{1-e_{uv}}, \quad e_{uv} \in \{0, 1\}. \quad (77)$$

To properly define a velocity $v(x, t)$, it should follow the continuity equation

$$\frac{\partial}{\partial t} p_t(e_{uv}) + \nabla \cdot (p v)_t(e_{uv}) = 0. \quad (78)$$

We use x and y to denote two states of e_{uv} ($p(e_{uv} = x) := p(x), p(e_{uv} = y) := p(y)$), then the divergence term is

$$\nabla \cdot (p v)(e_{uv} = x) = \sum_{y \neq x} [p_t(y) v_t(y \rightarrow x) - p_t(x) v_t(x \rightarrow y)]. \quad (79)$$

As we are working on a Bernoulli distribution, then the forward equations become

$$\begin{cases} \partial_t p(0) = p(1) v_t(1 \rightarrow 0) - p(0) v_t(0 \rightarrow 1), \\ \partial_t p(1) = p(0) v_t(0 \rightarrow 1) - p(1) v_t(1 \rightarrow 0). \end{cases} \quad (80)$$

Since $p_t(1) = [\mathbf{W}_t]_{uv}$, we have $\partial_t p(1) = [\dot{\mathbf{W}}_t]_{uv}$ and $\partial_t p(0) = -[\dot{\mathbf{W}}_t]_{uv}$. Hence

$$p(0) v_t(0 \rightarrow 1) - p(1) v_t(1 \rightarrow 0) = [\dot{\mathbf{W}}_t]_{uv}.$$

There are many solutions to the above equation. We chose a symmetric solution so that the transition of $e_{uv} \rightarrow 1 - e_{uv}$ with

$$v_t(0 \rightarrow 1) = \frac{[\dot{\mathbf{W}}_t]_{uv}}{1 - [\mathbf{W}_t]_{uv}}, \quad v_t(1 \rightarrow 0) = -\frac{[\dot{\mathbf{W}}_t]_{uv}}{[\mathbf{W}_t]_{uv}}.$$

Finally for concise, we can write it as a velocity field on states $e_{uv} \in \{0, 1\}$, note $1 - 2e_{uv}$ is +1 at $e_{uv} = 0$ and -1 at $e_{uv} = 1$. Thus, we have

$$v(e_{uv}, t) = (1 - 2e_{uv}) \frac{[\dot{\mathbf{W}}_t]_{uv}}{[\mathbf{W}_t]_{uv} (1 - [\mathbf{W}_t]_{uv})}, e_{uv} \in \{0, 1\},$$

which in matrix form gives

$$v_t(E_t | G_1, G_0) = (1 - 2E_t) \frac{\dot{\mathbf{W}}_t}{\mathbf{W}_t \circ (1 - \mathbf{W}_t)}. \quad (81)$$

Combine the node velocity and the edge velocity, we can now introduce the Discrete Bures-Wasserstein Flow matching algorithm, with the training and inference part respectively introduced in Algorithm 3 and Algorithm 4.

E Design space for Bures Wasserstein interpolation and velocity

In the introduction part, we have already compared different probability paths and how they are impacting the inference time sampling. While the Bures-Wasserstein flow path is shown to produce a better probability path for the model to learn, as we illustrated in Fig. 1a, we have to point out that linear interpolation and the corresponding probability path can still converge to the data distribution with sufficiently large flow steps. As if we conduct sampling with infinite flow steps during the later stage of flow, the samples are still able to arrive at the target distributions. A similar pattern exists in diffusion models when they are considered as a Monte-Carlo Markov Chain, and they need sufficiently large steps to converge. We emphasize that the convergence gap in Fig. 1c would be slowly recovered as the number of flow steps increases.

Algorithm 3: Discrete BWFlow Training**Input:** Ref. dist p_0 and dataset $\mathcal{D} \sim p_1$.**Output:** Trained model $f_\theta(G_t, t)$.

```

1 Initialize model  $f_\theta(G_t, t)$ ;
2 while  $f_\theta$  not converged do
  /* Sample Boundary Graphs */
3  Sample batched  $\{G_0\} \sim p_0, \{G_1\} \sim \mathcal{D}$ ;
  /* Prob.path Construction */
4  Sample  $t \sim \mathcal{U}(0, 1)$ ;
5  Calculate the BW interpolation to obtain
    $\mathbf{X}_t, \mathbf{W}_t$  via Eq. (11);
6  Sample  $G_t \sim p(\mathcal{G}_t | G_0, G_1)$  according
   to Eq. (13);
  /* Denoising - x-prediction */
7   $p_{1|t}^\theta(\cdot | G_t) \leftarrow f_\theta(G_t, t)$ ;
8  Loss calculation via Eq. (4);
9  optimizer.step();

```

Algorithm 4: Discrete BWFlow Sampling**Input:** Reference distribution p_0 , TrainedModel $f_\theta(G_t, t)$, Small time step dt ,**Output:** Generated Graphs $\{\hat{G}_1\}$.

```

1 Initialize samples  $\{\hat{G}_0\} \sim p_0$ ;
2 Initialize the model  $p_{1|t}^\theta(\cdot | G_t) \leftarrow f_\theta(G_t, t)$ 
  for  $t \leftarrow 0$  to  $1 - dt$  by  $dt$  do
  /* Denoising - x-prediction */
3  Predict  $\tilde{G}_1 \leftarrow p_{1|t}^\theta(\cdot | G_t)$ ;
  /* Velocity calculation */
4  Calculate  $v_\theta(\hat{G}_t | \hat{G}_0, \tilde{G}_1)$  via Eq. (14);
  /* Numerical Sampling */
5  Sample  $\hat{G}_{t+dt} \sim \hat{G}_t + v_\theta(\hat{G}_t)dt$ 

```

Given that different sampling algorithms can all bring the samples to the data distributions under certain conditions, we wish to understand the huge design space of Bures Wasserstein interpolation. We list the advantages and disadvantages in different techniques and discuss further when each techniques should be used.

In general, we consider two important steps to construct the flow matching for graph generation, specifically, *training* and *sampling*. In training, the main challenge is to obtain a valid real velocity $u(G_t)$ to be regressed to, so we listed a few strategies that can help us with that. In sampling, the challenge becomes how to reconstruct the probability path through the velocity estimated.

E.1 The training design

In general, the learning objective in flow matching depends on regressing the velocity term. There are several way to obtain the velocity.

1. **Exact velocity estimation.** Use Eq. (3) as the parameterization and learn $p_\theta(\mathcal{G}_1 | G_t)$
2. **Numerical approximation.** In the implementation of [50], the derivative is calculated through numerical approximation. To achieve better efficiency in calculating velocity, we simply consider a numerical estimation as in [50], where the velocity term is obtained as, $\dot{\mathbf{L}}_t = (\mathbf{L}_{t+\Delta t} - \mathbf{L}_t)/\Delta t$. Regressing on the numerical difference can provide an estimation for the velocity.
3. **AutoDiff.** In [11], the derivative of the probability path is evaluated through Pytorch AutoDiff. However, in practice we find this method unstable.

We summarized the training stage model parameterization in Table 4

	Continuous Flow Matching	Discrete Flow Matching
x-prediction	$v_t^\theta(G_t) = \frac{1}{1-t} [\tilde{G}_t^\theta(G_t) - G_t]$	$v_t^\theta(G_t) = \frac{1}{1-t} [p_{1 t}^\theta(\mathcal{G}_1 G_t) - \delta(\cdot, G_t)]$
Numerical Approximation	$v_t^\theta(G_t) \approx G_{t+dt} - G_t$	$v_t^\theta(G_t) \approx p(G_{t+dt} G_0, G_1) - p(G_t G_0, G_1)$
AutoDiff	$v_t^\theta(G_t) \approx \dot{G}_t$	Discrete velocity introduced in Eq. (14)

Table 4: The model parameterization for flow matching in training stage

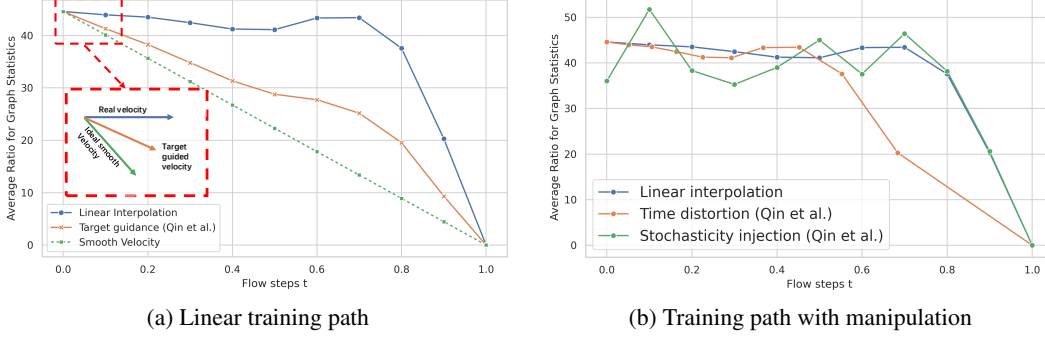


Figure 4: Techniques for manipulating probability path.

E.2 The sampling design

As we described in the Eq. (3), in our training framework, we actually train a denoised $p_\theta(\mathcal{G}_1 | \mathcal{G}_t)$. With such a parameterization and taking discrete flow matching as an example, the sampling can be done through one of the following design choices:

1. **Target-guided velocity sampling.** The velocity is designed as,

$$v_\theta(G_t) = \frac{1}{1-t} (p_\theta(\mathcal{G}_1 | G_t) - \delta(G_t, \cdot)).$$

This design directly moves the current point G_t towards the direction pointing to the predicted G_1 . The target-guided velocity is guaranteed to converge to the data distribution, but the interpolant might lie outside the valid graph domain.

2. **BW velocity sampling.** We use Eq. (14) to directly estimate the velocity and flow the Bures-Wasserstein probability path to generate new data points. This path is smooth in the sense of graph domain. However, this path requires more computational cost.
3. **Probability path reconstruction.** The third option is directly reconstructing the probability path, i.e., we first obtain an estimated point,

$$\tilde{G}_1 \sim p_\theta(\mathcal{G}_1 | G_t)$$

and then construct the data point at $t + dt$, which gives

$$G_{t+dt} \sim p(\mathcal{G}_t | \tilde{G}_1, G_0)$$

through Eq. (12). This is the most computationally costly method, which is obtained through the diffusion models. But this method also provide accurate probability path reconstruction.

In Section 4, we show BW velocity follows a path that minimizes the Wasserstein distance thus provides better performance, but sampling following linear velocity also provides convergence with much lower computational cost. So it is a trade-off to be considered in the real-world application.

	Continuous Flow Matching	Discrete Flow Matching
BW Velocity	Eq. (12)	Eq. (14)
Target-guided Velocity	$v_\theta(G_t) = \frac{1}{1-t} (\tilde{G}_1 - G_t)$	$v_\theta(G_t) = \frac{1}{1-t} (p_\theta(\mathcal{G}_1 G_t) - \delta(G_t, \cdot))$
Path Reconstruction	$G_{t+dt} \sim p(G_{t+dt} \tilde{G}_1, G_0)$	$G_{t+dt} \sim p(G_{t+dt} \tilde{G}_1, G_0)$

Table 5: Reconstructing probability path choices in flow matching during inference

F Discussion and limitations

F.1 The implicit manipulation of probability path

Though not explicitly mentioned, Qin et al. [45] makes huge efforts to manipulate the probability path for better velocity estimation by extensively searching the design space, and their finding aligns well

with the statement that the velocity should be smooth and consistently directing to the data points: 1) **Time distortion:** (The orange line in Fig. 4b) the polynomial distortion of training and sampling focus on the later stage of the probability trajectory, providing better velocity estimation in this area. This uneven sampling strategy is equivalent to pushing the probability path left to make it smooth. 2) **Target guidance:** (The orange line in Fig. 4a) the target guidance directly estimate the direction from a point along the path towards the termination graph, so that the manipulated probability could smoothly pointing to the data distribution. and 3) **Stochasticity injection:** (The green line in Fig. 4a) Stochasticity explores the points aside from the path, which avoid the path to be stuck in the platform area.

F.2 Potential extension to diffusion models

In order to extend the flow matching algorithms with diffusion models, one important thing is to convert the pair-conditioned probability path and velocity to single boundary conditions. For instance, the probability path in flow matching has the form $p(G_t | G_1, G_0)$ and the velocity follows $v(G_t | G_1, G_0)$. As suggested in [48, 8, 60], the discrete graph diffusion models require a velocity (which is equivalent to a ratio matrix) to perturb the data distribution conditioned on the data points, which we denote as $v(G_t | G_1)$. As long as the unilateral conditional velocity has a tractable form, one can first sample a G_1 and get G_t through iteratively doing to:

$$G_{t-dt} = G_t - v(G_t | G_1)dt$$

starting from G_1 . So that one can easily construct the probability path $p(G_t | G_1)$ to fit into the diffusion model framework. In practice, given that we know the explicit form of $v(G_t | G_1, G_{t'})$ (just replace G_0 in the expression), the unilateral conditional velocity can be obtained through taking the limitation,

$$v(G_t | G_1) = v(G_t | G_1, G_t) = \lim_{t' \rightarrow t} v(G_t | G_1, G_{t'}).$$

Both linear interpolation and our Bures-Wasserstein interpolation can achieve this easily. We just provide a discussion here and will leave this as future work as this paper does not focus on diffusion models but on flow matching models.

F.3 Permutation invariance

The Bures-Wasserstein distance between two graph distributions is not permutation invariant, and the minimal value is obtained through the graph alignment. So ideally, to achieve optimal transport, graph alignment and mini-batch matching could provide a better probability path. However, permutation invariance is not always a desired property since we only want to find a path that better transforms from the reference distribution to the data distributions. As an illustration, the widely used linear interpolation to construct graph flow [45] does not guarantee permutation invariance as well. And it is proved that, if the measurement is based on Wasserstein distance between two Gaussian distributions.

$$\begin{aligned} d_{\text{BW}}(\eta_{\mathcal{G}_0}, \eta_{\mathcal{G}_1}) &\leq d_{\text{Arithmetic}}(\eta_{\mathcal{G}_0}, \eta_{\mathcal{G}_1}) \\ \text{with } d_{\text{BW}}(\eta_{\mathcal{G}_0}, \eta_{\mathcal{G}_1}) &= \|\mathbf{X}_0 - \mathbf{X}_1\|_F^2 + \beta \text{trace} \left(\mathbf{L}_0^\dagger + \mathbf{L}_1^\dagger - 2 \left(\mathbf{L}_0^{\dagger/2} (\mathbf{P}^\top \mathbf{L}_1 \mathbf{P})^\dagger \mathbf{L}_0^{\dagger/2} \right)^{1/2} \right), \\ \text{and } d_{\text{Arithmetic}}(\eta_{\mathcal{G}_0}, \eta_{\mathcal{G}_1}) &= \|\mathbf{X}_0 - \mathbf{X}_1\|_F^2 + \beta \|\mathbf{L}_0 - \mathbf{P}^\top \mathbf{L}_1 \mathbf{P}\|^2, \forall \mathbf{P} \in \text{potential permutation set} \end{aligned} \quad (82)$$

G Related works

G.1 Diffusion and flow Models

Among contemporary generative models, diffusion [24] and flow models [33] have emerged as two compelling approaches for their superior performance in generating text and images. In particular, these generative models can be unified under the framework of stochastic interpolation [1], which consists of four procedures [34] as we introduced in Section 1. These contemporary generative models rely on constructing a probability path between data points of an easy-to-sample reference distribution and of the data distribution, and training a machine learning model to simulate the process [34]. So that one can sample from the reference (a.k.a source) distribution and iteratively

transform it to approximate data samples from the target distribution. Diffusion models construct the probability path with a unilateral path conditioned on the data distribution, where one start sampling from a data point X_1 and construct the path $p(\mathcal{X}_t | X_1)$. While flow models can condition on either two boundary conditions, $\{X_1, X_0\}$ or just one-side boundary condition X_1 .

Depending on the space that the algorithm operates on, both models can be categorized into continuous or discrete models. The continuous generative models assume the data distributions are themselves lying in continuous space (such as Gaussian) and build models, with examples in diffusion [24, 49, 58] and flow [33, 36]. The discrete generative models assume the data follows a discrete distribution, for instance categorical or Bernoulli distributions. Examples include discrete diffusion [8, 51] and discrete flow models [9, 19, 40].

Under the stochastic interpolation framework, the interpolation methods are commonly selected through optimal transport (OT) displacement interpolant [36, 1, 39]. Optimal transport is a classical topic in mathematics that was originally used in economics and operations research [57], and has now become a popular tool in generative models. OT aims for finding the best transport plan between two probability measures with the smallest associated transportation cost. It has been shown that generative models can be combined with technologies such as iterative matching [54] and mini batching [44] to approximate the OT cost, and get a significant boost in their performance in generative modeling.

G.2 Graph generation models

Thanks to the capability of graphs in representing complex relationships, graph generation [65, 35] has become an essential task in various fields such as protein design [26], drug discovery [6], and social network analysis [31]. The initial attempt at graph generation is formalized through autoregression. For instance, GraphRNN [62] organizes the node interactions into a series of connection events and conducts autoregressive prediction for generation. Later, one shot generation methods such as Variational Graph Auto-Encoder were proposed [30, 10].

Among various generative models, diffusion models and flow-based models have emerged as two compelling approaches for their ability to achieve state-of-the-art performance in graph generation tasks [41, 55, 18, 45, 25]. In the early stage, continuous diffusion models were first extended to the task of graph generation [41], where they just view the adjacency matrix as a special signal living on the $\mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ domain. However, these methods fail to capture the natural discreteness of graphs, and Vignac et al. [55] first brings discrete diffusion into graph generation. After that, more work [48, 60] starts to focus on designing better discrete diffusion models for graph generation.

On the other hands, with the development of flow matching techniques, a few works have been developed to utilize flow models for graph generation and they have achieved huge success. Eijkelboom et al. [18] utilizes variational flow matching to process categorical data and Qin et al. [45] developed discrete flow matching for graph generation tasks.

In parallel, there are a number of work that have managed to respect the intrinsic nature of graphs, such as global patterns. For instance, Jo et al. [28] brings a mixture of graph technique to enhance the performance by explicitly learning final graph structures; Yu and Zhan [63] mitigates exposure bias and reverse-start bias in graph generation; Hou et al. [25] improves graph generation through optimal transport flow matching techniques but they still assume the independence between nodes and edges and use hamming distance to measure the transport cost; and Li et al. [31] gives the large-scale attributed graph generation framework through batching edges.

However, there remain a core challenge: constructing the probability path p_t . Existing text and image generative models, operating either in the continuous [24, 49, 33, 36] or discrete [8, 51, 9, 19, 40] space, typically rely on linear interpolation between source and target distributions to construct the path. Graph generation models, including diffusion [41, 55, 23, 60, 48] and flow-based models [18, 45, 25], inherit this design by modeling every single node and edge independently and linearly build paths in the disjoint space. However, this approach is inefficient because it neglects the strong interactions and relational structure inherent in graphs, i.e., the significance of a node heavily depends on the configuration of its neighbors. While empirical success have been achieved via fine-grained searching on the training and sampling design [45] such as target guidance and time distortion, we argue that there remains a fundamental issue of the linear probability path construction, and these strategies only mitigate the problem by manipulating the probability path.

Table 6: Training and sampling time on each dataset. TG means using target-guided velocity; BW means using BW velocity.

Dataset	Min Nodes	Max Nodes	Training Time (h)	Graphs Sampled	Sampling Time (h)
Planar	64	64	45 (1.55x)	40	0.07(TG); 0.13 (BW)
Tree	64	64	10 (1.25x)	40	0.07(TG);0.14(BW)
SBM	44	187	74 (0.98x)	40	0.07(TG) 0.14(BW)
Moses	8	27	35 (0.76x)	25000	5(TG); 6(BW)
Guacamol	2	88	251 (1.8x)	10000	7(TG); 21(BW)
QM9	3	29	15	25000	5(TG) 6(BW)
GEOM	/	181	141	10000	7(TG) 14(BW)

Table 7: Best Configuration for Training and Sampling when producing Tables 1 to 3.

Dataset	Train		Sampling	
	Initial Distribution	Train Distortion	Sample Distortion	Sampling steps
Planar	Marginal	Identity	Identity	1000
Tree	Marginal	Polydec	Polydec	1000
SBM	Absorbing	Identity	Identity	1000
MOSES	Marginal	Identity	Identity	500
GUACAMOL	Marginal	Identity	Identity	500
QM9	Marginal	Identity	Identity	500
GEOM	Marginal	Identity	Identity	500

H Comparison with other interpolation methods

In the experimental part, we compare our methods with arithmetic (linear) interpolation, geometric interpolation and harmonic interpolation. We state the equation of them respectively as follows.

We consider the boundary graph G_0 and G_1 with $\mathbf{X}_0, \mathbf{X}_1 \in \mathbb{R}^{|\mathcal{V}| \times d}$ and $\mathbf{W}_0, \mathbf{W}_1 \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$. Let $t \in [0, 1]$, we fixed the feature interpolation as,

$$\mathbf{X}_t = (1 - t) \mathbf{X}_0 + t \mathbf{X}_1,$$

the graph structure interpolation can be expressed as,

Linear interpolation:

$$\mathbf{W}_t = (1 - t) \mathbf{W}_0 + t \mathbf{W}_1.$$

Geometric interpolation:

$$\mathbf{W}_t = \mathbf{W}_0^{1/2} (\mathbf{W}_0^{-1/2} \mathbf{W}_1 \mathbf{W}_0^{-1/2})^t \mathbf{W}_0^{1/2},$$

Harmonic interpolation:

$$\mathbf{W}_t = \left((1 - t) \mathbf{W}_0^{-1} + t \mathbf{W}_1^{-1} \right)^{-1}.$$

Each interpolation methods actually handle each special manifold assumption, which should be designed under a comprehensive understanding of the task. In our experimental part, we conduct intensive analysis on the impact of interpolation methods to the graph generation quality.

I Additional experiment results

I.1 Experiment setups and computational cost

The training and sampling computation time are provided in Table 6. The experiments were run on a single NVIDIA A100-SXM4-80GB GPU. The hyperparameter configuration in producing Tables 1 to 3 is reported in Table 7.

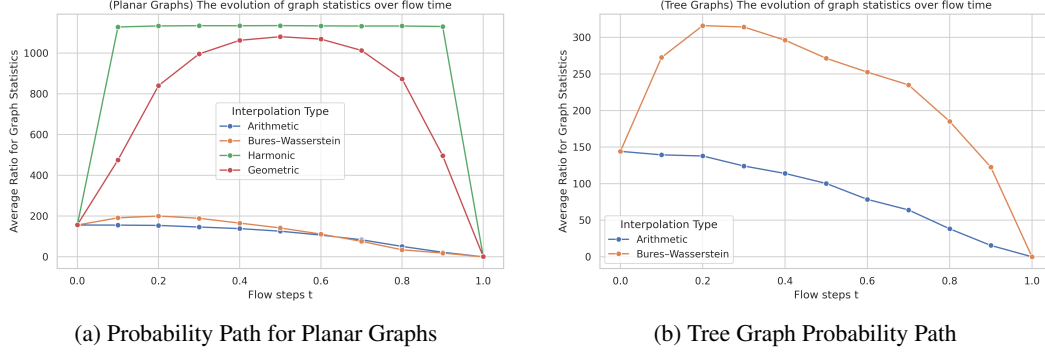


Figure 5: BW probability paths for planar and tree graphs.

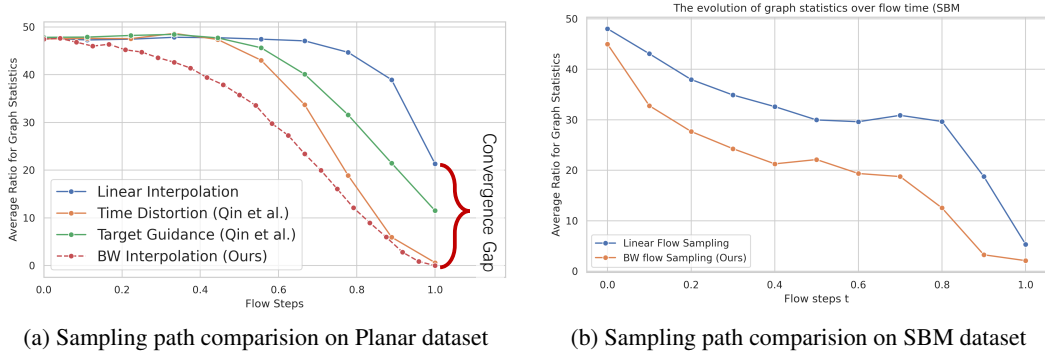


Figure 6: The probability path reconstruction in the sampling stage on a) Planar graphs and b) SBM graphs.

I.2 Additional results for the training paths

Fig. 5 gives the training probability path construction for planar graphs and tree graphs. While planar graphs have a similar pattern as the SBM datasets as in Fig. 3a, the probability path constructed for tree graphs does not follow a similar pattern. We attribute this to the different geometry of tree graphs that reside in hyperbolic space [61].

I.3 More experiments on plain graph generations

Additional results for sampling paths. We then give the sampling path construction in Fig. 6. To better illustrate the advantage of BWFlow, we fix the sampling steps to be as small as 50. It is clear that in planar and SBM dataset, the BW velocity can still provide a smooth probability and stable convergence towards the data distribution. While the linear velocity does not give a good probability path and fails to converge to the optimal value, especially when the sampling size is small.

The maximum mean discrepancy (MMD) of four graph statistics between the set of generated graphs and the test set is measured, including degree (Deg.), clustering coefficient (Clus.), count of orbits with 4 nodes (Orbit), the eigenvalues of the graph Laplacian (Spec.), wavelet ratio (Wavelet.). To verify that the model learns to generate graphs with valid topology, we give the percentage of valid, unique, and novel (V.U.N.) graphs for where a valid graph satisfies the corresponding property of each dataset (Planar, Tree, SBM, etc.).

Full results for plain graph generation. Table 8 gives the full results with other generative models aside from the diffusion and flow models. Table 9 gives the results on smaller datasets, i.e., comm20

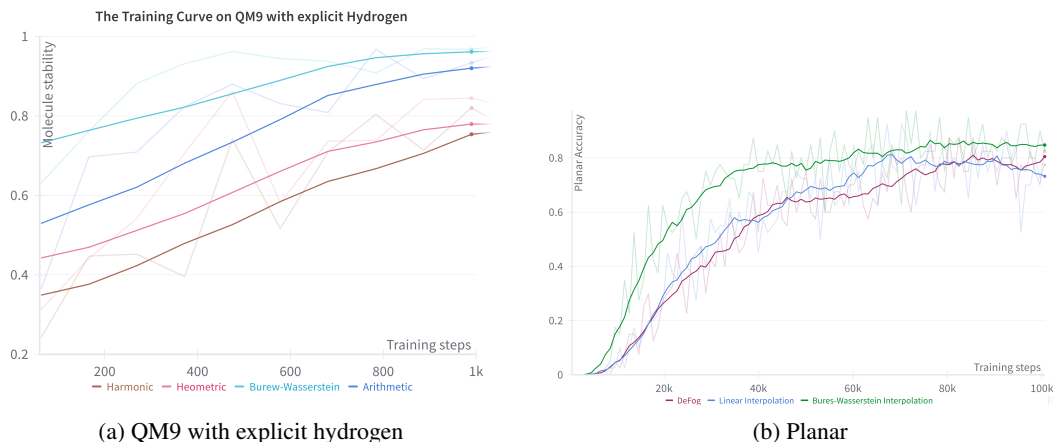


Figure 7: Training curves on QM9 and planar datasets with explicit hydrogen.

I.4 Full table for the synthetic graph generation

In Table 10, we illustrate the numerical results for comparing the interpolation methods in plain graph generation without node features. It is clear that BWFlow outperforms other methods in planar and SBM graphs. But the performance was not good in tree graph generations.

I.5 3D Molecule Generation: QM9 without explicit hydrogen

In Table 11 we report the results of QM9 without explicit hydrogen. This task is relatively easy compared to the generation task with explicit hydrogen, and both Midi and our BWFlow have achieved near-saturated performance with validity near to 100%.

I.6 Convergence Analysis

Fig. 7 are the training convergence analysis on Planar and QM9 dataset, showing that BWFlow provides a fast convergence speed than others.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: The abstract summarize the contributions and the novelty the method.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: In the final section we provided a discussion of the limitations.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: Each claim comes with a remark of the assumptions made, and proof are made in appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: We provide experimental details in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: Code can be provided upon request, and will be released when the paper moves to the following stages.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Experimental details is provided in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We provided the standard deviation together with the mean of the experiment results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide the computation architecture details and the training time in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: We follow every detail of the guidelines and the code of conduct.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: The core contribution of the paper is the theoretical results and a novel method. The societal impacts are not applicable here.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We do not propose new datasets, and the model do not pose safeguard risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We explicitly credit our relevant assets used.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. **New assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: No newly released assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and research with human subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: We do not conduct our research on human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: Not applicable.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. **Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: We do not utilize LLMs for developing the core components of the research.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.

Table 8: Graph Generation Performance on Synthetic Graphs. Results are obtained through tuning the probability path manipulation techniques. The remaining values are obtained from Qin et al. [45].

Planar Dataset					
Model	Ratio ↓	Valid ↑	Unique ↑	Novel ↑	V.U.N. ↑
Train set	1.0	100	100	0.0	0.0
GRAN [32]	2.0	97.5	85.0	2.5	0.0
SPECTRE [38]	3.0	25.0	100	100	25.0
DiGress [55]	5.1	77.5	100	100	77.5
EDGE [13]	431.4	0.0	100	100	0.0
BwR [15]	251.9	0.0	100	100	0.0
BiGG [14]	16.0	62.5	85.0	42.5	5.0
GraphGen [20]	210.3	7.5	100	100	7.5
HSpectre (one-shot) [4]	1.7	67.5	100	100	67.5
HSpectre [4]	2.1	95.0	100	100	95.0
GruM [28] 1.8	—	—	—	90.0	—
CatFlow [18]	—	—	—	—	80.0
DisCo [60]	—	83.6 ±2.1	100.0 ±0.0	100.0 ±0.0	83.6 ±2.1
Cometh - PC [48]	—	99.5 ±0.9	100.0 ±0.0	100.0 ±0.0	99.5 ±0.9
DeFoG	1.6 ±0.4	99.5 ±1.0	100.0 ±0.0	100.0 ±0.0	99.5 ±1.0
BWFlow	1.3 ±0.4	97.5 ±2.5	100.0 ±0.0	100.0 ±0.0	97.5 ±2.5

Tree Dataset					
Train set	1.0	100	100	0.0	0.0
GRAN [32]	607.0	0.0	100	100	0.0
DiGress [55]	1.6	90.0	100	100	90.0
EDGE [13]	850.7	0.0	7.5	100	0.0
BwR [15]	11.4	0.0	100	100	0.0
BiGG [14]	5.2	100	87.5	50.0	75.0
GraphGen [20]	33.2	95.0	100	100	95.0
HSpectre (one-shot) [4]	2.1	82.5	100	100	82.5
HSpectre [4]	4.0	100	100	100	100
DeFoG	1.6 ±0.4	96.5 ±2.6	100.0 ±0.0	100.0 ±0.0	96.5 ±2.6
BWFlow	1.4 ±0.3	75.5 ±2.4	100.0 ±0	100.0 ±0	75.5 ±2.4

Stochastic Block Model ($n_{\max} = 187, n_{\text{avg}} = 104$)					
Model	Ratio ↓	Valid ↑	Unique ↑	Novel ↑	V.U.N. ↑
Training set	1.0	85.9	100	0.0	0.0
GraphRNN [62]	14.7	5.0	100	100	5.0
GRAN [32]	9.7	25.0	100	100	25.0
SPECTRE [38]	2.2	52.5	100	100	52.5
DiGress [55]	1.7	60.0	100	100	60.0
EDGE [13]	51.4	0.0	100	100	0.0
BwR [15]	38.6	7.5	100	100	7.5
BiGG [14]	11.9	10.0	100	100	10.0
GraphGen [20]	48.8	5.0	100	100	5.0
HSpectre (one-shot) [4]	10.5	75.0	100	100	75.0
HSpectre [4]	10.2	45.0	100	100	45.0
GruM [28]	1.1	—	—	—	85.0
CatFlow [18]	—	—	—	—	85.0
DisCo [60]	—	66.2 ±1.4	100.0 ±0.0	100.0 ±0.0	66.2 ±1.4
Cometh [48]	—	75.0 ±3.7	100.0 ±0.0	100.0 ±0.0	75.0 ±3.7
DeFoG	4.9 ±1.3	90.0 ±5.1	100.0 ±0.0	100.0 ±0.0	90.0 ±5.1
BWFlow	3.8 ±0.9	90.5 ±4.0	100.0 ±0.0	100.0 ±0.0	90.5 ±4.0

FM type	Interpolation	comm-20		
		Deg.	Clus.	Orbit.
Discrete	Arithmetic	0.071	0.115	0.037
	Harmonic	0.011	0.036	0.027
	Geometric	0.047	0.083	0.02
	BW	0.009	0.013	0.017

Table 9: Quantitative experimental results on COMM20 (smaller dataset).

Dataset	FM type	Interpolation	V.U.N metrics			Spectral Metrics					
			Novelty	Uniqueness	Validity	Orbit	Spec	Clustering	Degree	Wavelet	Avg. Ratio
Planar	Discrete	DeFog	100	100	78.25	8.98	1.45	2.09	2.65	2.38	3.51
		Arithmetic	100	100	73.25	10.83	1.33	1.74	2.24	2.39	3.70
		harmonic	100	100	0.00	4519.63	2.57	17.01	25.10	42.41	921.35
		Geometric	100	100	23.25	655.66	1.61	10.17	13.25	6.68	137.47
		BW	100	100	84.75	5.14	1.27	1.69	1.78	2.02	2.38
Tree	Discrete	Defog	100	100	61.53	/	1.17	/	1.27	1.51	1.32
		Arithmetic	100	100	56.91	/	1.16	/	1.04	1.45	1.22
		harmonic	100	100	0.53	/	2.32	/	1.93	3.31	2.52
		Geometric	100	100	48.38	/	1.62	/	2.13	2.10	1.94
		BW	100	100	51.45	/	1.58	/	2.56	2.13	2.09
SBM	Discrete	Arithmetic	100	100	44.63	2.57	1.40	1.46	15.55	7.88	5.77
		harmonic	100	100	9.73	3.10	10.23	1.59	172.10	103.04	58.10
		Geometric	100	100	0.0	3.11	4.45	1.80	150.41	60.60	54.65
		BW	100	100	58.70	2.03	1.50	1.50	9.04	8.41	4.51

Table 10: Ablation study on interpolation methods when probability path manipulation techniques are all disabled. The clustering and orbit ratios in tree graphs are omitted, given that in the training set, the corresponding statistics are 0. The results go over exponential moving average (decay 0.999) for the last 5 checkpoints. The table is produced with Marginal boundary distributions, without time distortion.

Dataset	Interpolation	Metrics					
		μ	V.U.N(%)	Connected(%)	Charges(10^{-2})	Atom(10^{-2})	Angles($^{\circ}$)
QM9 (w/o h)	MiDi	1.00	98.0	100.0	0.4	5.1	1.49
	Linear Flow	1.60	79.33	52.3	0.7	14.0	8.77
	BWFlow	1.02	99.8	100.0	0.4	4.8	1.53

Table 11: Quantitative experimental results on QM9 datasets without explicit hydrogen in 3D molecule generation.

Table 12: Graph generation performance on the synthetic datasets: Planar, Tree, and SBM. Given that the synthetic datasets are usually unstable in evaluation, we applied an exponential moving average to stabilize the results and sample 5 times (each run generates 40 graphs) to calculate the mean and standard deviation. The experiment settings are in Table 7.

Model	Class	Planar		Tree		SBM	
		V.U.N. \uparrow	A.Ratio \downarrow	V.U.N. \uparrow	A.Ratio \downarrow	V.U.N. \uparrow	A.Ratio \downarrow
Train set	—	100	1.0	100	1.0	85.9	1.0
DiGress [55]	Diffusion	77.5	5.1	90.0	1.6	60.0	<u>1.7</u>
EDGE [13]	Diffusion	0.0	431.4	0.0	850.7	0.0	51.4
HSpectre [4]	Diffusion	<u>95.0</u>	2.1	100.0	4.0	75.0	10.5
GruM [28]	Diffusion	90.0	<u>1.8</u>	—	—	85.0	1.1
CatFlow [18]	Flow	80.0	—	—	—	85.0	—
DisCo [60]	Diffusion	83.6 \pm 2.1	—	—	—	66.2 \pm 1.4	—
Cometh [48]	Diffusion	99.5 \pm 0.9	—	—	—	75.0 \pm 3.7	—
DeFoG [45]	Flow	99.5 \pm 1.0	1.6 \pm 0.4	96.5 \pm 2.6	1.6 \pm 0.4	90.0 \pm 5.1	4.9 \pm 1.3
BWFlow	Flow	<u>97.5</u> \pm 2.5	1.3 \pm 0.4	75.5 \pm 2.4	1.4 \pm 0.3	90.5 \pm 4.0	3.8 \pm 0.9
DeFoG (EMA)	Flow	77.5 \pm 8.37	3.5 \pm 1.7	73.1 \pm 11.4	1.50 \pm 0.3	85.0 \pm 7.1	3.7 \pm 0.9
BWFlow (EMA)	Flow	84.8 \pm 6.44	2.4 \pm 0.9	68.1 \pm 12.2	1.24 \pm 0.2	83.5 \pm 6.0	2.4 \pm 0.6

Table 13: Comparison of interpolation methods on 3D Molecule generation with explicit hydrogen in QM9 dataset.

Flow Type	Interpolation	Metrics							
		μ	V.U.N(%)	Mol Stable	Atom Stable	Connected(%)	Charge(10^{-2})	Atom(10^{-2})	Angles($^{\circ}$)
Discrete	MiDi	1.01	93.13	93.98	99.60	99.21	0.2	3.7	2.21
	Arithmetic	1.01	87.53	88.45	99.13	99.09	0.4	4.2	2.72
	harmonic	1.01	94.91	94.54	99.65	99.03	0.6	6.4	2.21
	Geometric	1.01	91.26	91.29	99.42	98.42	0.1	4.4	3.63
	BW	1.01	96.45	97.84	99.84	99.24	0.1	2.3	1.96
Continuous	Arith	2.15	25.45	10.23	76.85	28.82	0.7	5.6	14.47
	harmonic	1.01	11.38	11.64	73.48	99.65	1.2	17.2	15.04
	Geometric	1.00	42.07	46.08	91.13	99.87	1.0	12.7	8.03
	BW	1.02	62.02	61.76	95.99	97.72	0.6	8.7	7.80

* Clearly, continuous flow matching models are not as comparative as discrete flow matching models.

Table 14: Large molecule generation results. Only iterative denoising-based methods are reported.

Model	Guacamol					MOSES						
	Val. \uparrow	V.U. \uparrow	V.U.N. \uparrow	KL div \uparrow	FCD \uparrow	Val. \uparrow	Unique. \uparrow	Novelty \uparrow	Filters \uparrow	FCD \downarrow	SNN \uparrow	Scaf \uparrow
Training set	100.0	100.0	0.0	99.9	92.8	100.0	100.0	0.0	100.0	0.01	0.64	99.1
DiGress [55]	85.2	85.2	85.1	92.9	68.0	85.7	100.0	95.0	97.1	1.19	0.52	14.8
DisCo [60]	86.6	86.6	86.5	92.6	59.7	88.3	100.0	97.7	95.6	1.44	0.50	15.1
Cometh [48]	98.9	98.9	97.6	96.7	72.7	90.5	99.9	92.6	99.1	1.27	0.54	16.0
DeFoG [45]	99.0	99.0	97.9	97.7	73.8	92.8	99.9	92.1	98.9	1.95	0.55	14.4
BWFlow (Ours)	98.8	98.9	97.4	/	69.2	92.0	100.0	94.5	98.4	1.32	0.56	15.3

Table 15: Large molecule generation results. Only comparing the representative diffusion and flow models. B.E. is the scenario that only considers binary edge types. The results are almost saturated, thus not very informative.

Model	Guacamol			MOSES		
	Val. \uparrow	V.U. \uparrow	V.U.N. \uparrow	Val. \uparrow	Unique. \uparrow	Novelty \uparrow
Digress (B.E.)	96.0	98.9	97.4	96.1	100	100
Defog (B.E.)	98.4	98.4	97.9	99.3	100	100
BWFlow (B.E.)	98.0	98.0	97.7	99.6	100	100