

GRAPHCG: UNSUPERVISED DISCOVERY OF STEERABLE FACTORS IN GRAPHS

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep generative models have been widely developed for graph data such as molecular graphs and point clouds. Yet, much less investigation has been carried out on understanding the learned latent space of deep graph generative models. Such understandings can open up a unified perspective and provide guidelines for essential tasks like controllable generation. To this end, this work develops a method called GraphCG for unsupervised discovery of steerable factors in latent space of deep graph generative models. We first examine the representation space of the recent deep generative models trained for graph data, and observe that the learned representation space is not perfectly disentangled. Thus, our method is designed for discovering steerable factors of graph data in a model-agnostic and task-agnostic manner. Specifically, GraphCG learns the semantic-rich directions via maximizing the corresponding mutual information, where the edited graph along the same direction will possess certain steerable factors. We conduct experiments on two types of graph data, molecular graphs and point clouds. Both the quantitative and qualitative results show the effectiveness of GraphCG for discovering steerable factors.

1 INTRODUCTION

Graph is a general format for many real-world data. For instance, molecules can be treated as graphs [10, 14] where the chemical atoms and bonds correspond to the topological nodes and edges respectively. Processing point clouds as graphs is also a popular strategy [53, 58], where points are viewed as nodes and edges are built among the nearest neighbors. Many existing works on deep generative models (DGMs) focus on modeling the graph data and improving the synthesis quality. However, understanding DGMs on graph and their learned representations has been much less explored, which may hinder the development of important applications like the controllable generation (also referred to as the data editing) and the discovery of interpretable data structure.

The graph controllable generation task refers to modifying the steerable factors of graph so as to obtain graphs with desired properties [9, 43]. This is an important task in many applications, but traditional methods (*e.g.*, manual editing) possess certain inherent limitations under certain circumstances. A classic example is molecule editing, which aims at modifying the substructures of molecules [38] and can relate to some key tactics in drug discovery like functional group change [13] and scaffold hopping [2, 23]. This is a routine task in pharmaceutical companies, yet, relying on domain experts for manual editing can be subjective or biased [9, 15]. Different from previous works, this paper aims to explore the unsupervised graph editing with DGMs. It can act as a good complementary module to conventional methods and bring many crucial benefits: (1) It enables the efficient graph editing in the large-scale setting. (2) It alleviates the requirements for extensive domain knowledge for factor change labeling. (3) It provides another perspective for editing preference, which reduces biases from the domain experts.

One core property relevant to the general unsupervised data editing using DGMs is the disentanglement. While there does not exist a widely-accepted definition of disentanglement, the key intuition [36] is that a disentangled representation should separate the distinct, informative, and steerable factors of variations in the data. Thus, the controllable generation task would become trivial with the disentangled DGMs as the backbone. Such a disentanglement assumption has been widely used in generative modeling on the image data, *e.g.*, β -VAE [19] learns disentangled representation by forcing the representation to be close to an isotropic unit Gaussian. However, it may introduce extra constraints on the formulations and expressiveness of DGMs [11, 19, 47, 60].

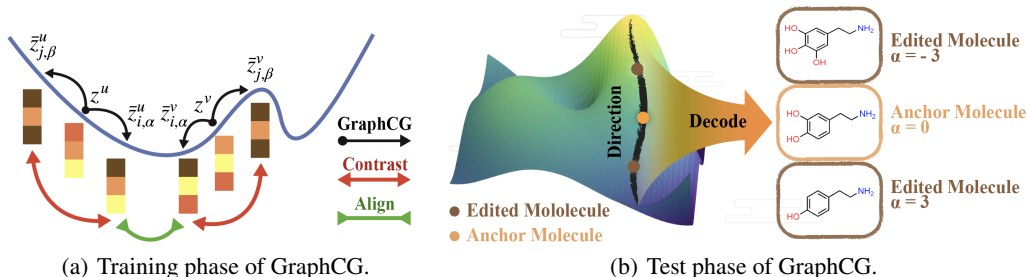


Figure 1: (a) The training phase. Given two latent codes z^u and z^v , we edit the four latent representations along i -th and j -th direction with step size α and β respectively. The goal of GraphCG is to align the positive pair ($\bar{z}_{i,\alpha}^u$ and $\bar{z}_{i,\alpha}^v$), and contrast them with $\bar{z}_{j,\beta}^u$ and $\bar{z}_{j,\beta}^v$ respectively. (b) The test phase. We will first sample an anchor molecule, and adopt the learned directions in the training phase for editing. With step size $\alpha \in [-3, 3]$, we can generate a sequence of molecules. Specifically, after decoding, there is a functional group change shown up: the number of hydroxyl groups decreases along the sequence in the decoded molecules.

For graph data, one crucial question remains: *is the latent representation space learned from DGMs on graph data disentangled?* In image generation, prior work [36] shows that without inductive bias, the latent space learned by VAEs is not guaranteed to be disentangled. However, the disentanglement property of graph DGMs is much less explored. In Sec. 3, we first study the latent space of DGMs on two typical graph data (molecular graphs and point clouds), and empirically illustrate that the learned space is not perfectly disentangled. This observation then raises the second question: *Given a pretrained DGM with not perfectly disentangled latent space, is there a flexible framework enabling the graph controllable generation in an unsupervised manner?* To tackle this, we propose a model-agnostic and task-agnostic framework called GraphCG for unsupervised graph controllable generation. GraphCG has two main phases, as illustrated in Fig. 1. During the training phase (Fig. 1(a)), GraphCG starts with the assumption that the steerable directions can be learned by maximizing the mutual information (MI) among the semantic directions. We formulate GraphCG with an energy-based model (EBM), which provides a large family of solutions. Then during the test phase, with the learned semantic directions, we can carry out the editing task by moving along the direction with certain step sizes. As the example illustrated in Fig. 1(b), the molecular structure (hydroxyl group) changes consistently along the editing sequence. For evaluation, we visually verify the learned semantic directions on both types of graph data. Further for the molecular graphs, we propose a novel evaluation metric called sequence monotonic ratio (SMR) to measure the output sequences.

We summarize our contributions as follows: (1) We conduct an empirical study on the disentanglement property of three pretrained deep generative models (DGMs) on two types of graph data, molecular graphs and point clouds. We find that the latent space of these pretrained graph DGMs is not perfectly disentangled. (2) We propose a model-agnostic and task-agnostic method called GraphCG for the unsupervised graph controllable generation. GraphCG aims at learning the semantic directions by maximizing their corresponding mutual information, and its outputs are sequences of graphs. (3) We qualitatively evaluate the proposed methods on two types of graph data, molecular graphs and point clouds. Besides, the quantitative results further show the clear improvement over the baselines.

Related work. Recent works leverage the DGMs for various controllable generation tasks [5, 61], where the inherent assumption is that the learned latent representations encode rich semantics, and thus traversal in the latent space can help steer factors of data [17, 26, 52]. Among them, one research direction [40, 52] is using supervised signals to learn the semantic-rich directions, and most works on editing the graph data focus on the supervised setting [27, 57, 64]. However, these approaches can not be applied to many realistic scenarios where extracting the supervised labels is difficult. Another research line [17, 45, 51] considers discovering the latent semantics in an unsupervised manner, but most unsupervised methods are designed to be either model-specific or task-specific, making them not directly applicable to the graph data. More comprehensive discussion is in Appendix B.

2 BACKGROUND AND PROBLEM FORMULATION

Graph and deep generative models (DGMs). Each graph data (including nodes and edges) is denoted as $x \in \mathcal{X}$, where \mathcal{X} is the data space, and DGMs learn the data distribution, *i.e.*, $p(x)$. Our

proposed graph editing method (GraphCG) is model-agnostic, so we briefly introduce the mainstream DGMs for graph data as below. Variational auto-encoder (VAE) [19, 31] measures a variational lower bound of $p(\mathbf{x})$ by introducing a proposal distribution; flow-based model [8, 46] constructs reversible encoding functions such that the data distribution can be deterministically mapped to a prior distribution. Note that these mainstream DGMs, either explicitly or implicitly, contain an encoder ($f(\cdot)$) and a decoder ($g(\cdot)$) parameterized by neural networks as:

$$\mathbf{z} = f(\mathbf{x}), \quad \mathbf{x}' = g(\mathbf{z}), \quad (1)$$

where $\mathbf{z} \in \mathcal{Z}$ is the latent representation, \mathcal{Z} is the latent space, and \mathbf{x}' is the reconstructed output graph. Since in the literature [51, 52], people also call latent representations as latent codes or latent vectors, in what follows, we will use these terms interchangeably.

Semantic direction and step size. In the latent space \mathcal{Z} , we assume there exist D semantically meaningful direction vectors, *i.e.*, \mathbf{d}_i with $i \in \{0, 1, \dots, D-1\}$. There is also a scalar variable, step size α , which controls the degree to edit the sampled data with desired steerable factors (as will be introduced below), and we follow the prior work [51] on taking $\alpha \in [-3, 3]$. Each direction corresponds to one or multiple factors, such that by editing the latent vector \mathbf{z} with direction \mathbf{d}_i and step size α , the reconstructed graph will be augmented with the desired factors, leading to certain structural changes.

Steerable factors. The steerable factors are attributes of DGMs, which usually refer to the semantic information of data that we can explicitly discover from the pretrained DGMs. In this work, we focus on the steerable factors of graph data, which are data- and task-specific. Yet, there is one category of factors that is commonly shared among all the graph data: the **structure** information. Concretely, these steerable factors can be the functional groups or fragments in molecular graphs and shapes or sizes in point clouds. In Appendix C, we provide a detailed description of these steerable factors.

Problem formulation: graph controllable generation. Given a pretrained DGM (*i.e.*, the encoder and decoder are fixed), our goal is to learn the most semantically rich directions (\mathbf{d}_i) in the latent space \mathcal{Z} . Then for each latent code \mathbf{z} , with the i -th semantic direction and a step size α , we can get an edited latent vector $\bar{\mathbf{z}}_{i,\alpha}$ and edited data $\bar{\mathbf{x}}'$, as:

$$\mathbf{z} = f(\mathbf{x}), \quad \bar{\mathbf{z}}_{i,\alpha} = h(\mathbf{z}, \mathbf{d}_i, \alpha), \quad \bar{\mathbf{x}}' = g(\bar{\mathbf{z}}_{i,\alpha}), \quad (2)$$

where \mathbf{d}_i and $h(\cdot)$ are the edit direction and edit function that we want to learn. We expect that $\bar{\mathbf{z}}_{i,\alpha}$ can inherently possess certain steerable factors, which can then be reflected in the graph structure of $\bar{\mathbf{x}}'$.

Energy-based model (EBM). EBM is a flexible framework for distribution modeling:

$$p(\mathbf{x}) = \frac{\exp(-E(\mathbf{x}))}{A} = \frac{\exp(-E(\mathbf{x}))}{\int_{\mathbf{x}} \exp(-E(\mathbf{x})) d\mathbf{x}}, \quad (3)$$

where $E(\cdot)$ is the energy function and A is the partition function. In EBM, the bottleneck is the estimation of partition function A : it is commonly intractable due to the high cardinality of \mathcal{X} . Various methods have been proposed to handle this issue, including but not limited to contrastive divergence [20], noise-contrastive estimation [4, 16], and score matching [24, 54, 56].

3 DISENTANGLEMENT OF LATENT REPRESENTATION

In this section, we quantify the degree of disentanglement of the existing DGMs for graph data. In specific, we adopt six disentanglement measures, and the observed low disentanglement scores show that compared to DGMs for image data, such as StyleGANs [28, 29], the attributes in latent space of graph DGMs are not well disentangled.

The key intuition [36] behind disentanglement is that a disentangled representation space should separate the distinct, informative, and steerable factors of variations in the data. In other words, each latent dimension of the disentangled representation corresponds to one or multiple factors. Therefore, the change of the disentangled dimension can lead to the consistent change in the corresponding factors of the data. This good property has become a foundational assumption in many existing controllable generation methods [17, 51, 52].

Is the latent space of graph DGMs disentangled? In image generation, [36] shows that without inductive bias, the representation learned by VAEs is not perfectly disentangled. To verify if this claim is also valid for the mainstream DGMs on graphs, we conduct the following experiment.

Table 1: The six disentanglement metrics on three pretrained DGMs and two graph types. All measures range from 0 to 1, and higher scores mean more disentangled representation.

Graph Type	Model	Dataset	BetaVAE \uparrow	FactorVAE \uparrow	MIG \uparrow	DCI \uparrow	Modularity \uparrow	SAP \uparrow
Molecular Graph	MoFlow	ZINC250K	0.260	0.175	0.031	0.953	0.620	0.009
	HierVAE	ChEMBL	0.178	0.165	0.022	0.114	0.606	0.026
Point Cloud	PointFlow	Airplane	0.022	0.025	0.029	0.160	0.745	0.022

Steerable factors and experiments on disentanglement measure. There have been a series of works exploring the disentanglement of the latent space in DGMs, and here we take six widely-used ones: BetaVAE [19], FactorVAE [30], MIG [6], DCI [11], Modularity [47], and SAP [32]. Each measure has its own bias, and we put a detailed comparison in Appendix C. Meanwhile, they all share the same high-level idea: given the latent representation from a pretrained DGM, they are proposed to measure how predictive it is to certain steerable factors.

To adapt them to our setting, first we need to extract the steerable factors in graph, which requires the domain knowledge. For instance, in molecular graphs, we can extract some special substructures named fragments or functional groups. These substructures can be treated as steerable factors since they are the key components of the molecules and are closely related to certain molecular properties [50]. We use RDKit [33] to extract 9 most distinguishable fragments as steerable factors for disentanglement measurement. For point clouds, we use PCL tool [48] to extract 75 VFH descriptors [49] as steerable factors, which depicts the geometries and viewpoints accordingly.

Then for measuring the disentanglement, we consider six metrics on two data types with three backbone models. All the metric values range from 0 to 1, and the higher the value, the more disentangled the DGM is. According to Table 1, We can observe that the most of disentanglement scores are quite low, except the DCI [11] on MoFlow. Thus, we can draw the conclusion that generally these graph DGMs are not perfectly disentangled. More details of this experiment (the steerable factors on two data types and six disentanglement metrics) can be found in Appendix C.

4 OUR METHOD

The analysis in Sec. 3 naturally raises the next research question: *given a not well-disentangled representation space, is there a flexible way to do the graph data editing?* The answer is positive. We propose GraphCG, a flexible model-agnostic and task-agnostic framework to learn the semantic directions in an unsupervised manner. It starts with the assumption that the latent representations edited with the same semantic direction and step size should possess similar information (corresponding to the factors) to certain degree, thus by maximizing the mutual information them, we can learn the most semantic-rich directions. Then we formulate this editing task as a density estimation problem with the energy-based model (EBM). As introduced in Sec. 2, EBM covers a broad range of solutions, and we further propose GraphCG-NCE by adopting the noise-contrastive estimation (NCE).

4.1 GRAPH CONTROLLABLE GENERATION WITH MUTUAL INFORMATION

The mutual information (MI) measures the non-linear dependency between variables. To adapt it to our setting, we set the editing condition as containing both the semantic directions and step sizes, and we assume that maximizing the MI between different conditions can maximize the shared information within each condition. The pipeline is as follows.

We first sample two codes in the latent space, \mathbf{z}^u and \mathbf{z}^v . Then we pick up the i -th semantic direction and one step size α to obtain the edited latent points in \mathcal{Z} . The corresponding points are as

$$\bar{\mathbf{z}}_{i,\alpha}^u = h(\mathbf{z}^u, \mathbf{d}_i, \alpha), \quad \bar{\mathbf{z}}_{i,\alpha}^v = h(\mathbf{z}^v, \mathbf{d}_i, \alpha). \quad (4)$$

Under our assumption, we expect that these two edited points share certain information with respect to the steerable factors. Thus, we want to maximize the MI between $\bar{\mathbf{z}}_{i,\alpha}^u$ and $\bar{\mathbf{z}}_{i,\alpha}^v$. Since the MI is intractable to compute, we adopt the EBM lower bound [35] as:

$$\mathcal{L}_{\text{MI}}(\bar{\mathbf{z}}_{i,\alpha}^u, \bar{\mathbf{z}}_{i,\alpha}^v) = \frac{1}{2} \mathbb{E}_{p(\bar{\mathbf{z}}_{i,\alpha}^u, \bar{\mathbf{z}}_{i,\alpha}^v)} [\log p(\bar{\mathbf{z}}_{i,\alpha}^u | \bar{\mathbf{z}}_{i,\alpha}^v) + \log p(\bar{\mathbf{z}}_{i,\alpha}^v | \bar{\mathbf{z}}_{i,\alpha}^u)]. \quad (5)$$

The detailed derivation is in Appendix D. Till this step, we have transformed the graph data editing task into the summation of two conditional log-likelihoods estimation problem.

4.2 GRAPHCG WITH ENERGY-BASED MODEL

Following Eq. (5), maximizing the MI between $I(\bar{z}_{i,\alpha}^u; \bar{z}_{i,\alpha}^v)$ is equivalent to estimating the summation of two conditional log-likelihoods. We then model them using two conditional EBMs. Because these two views are in the mirroring direction, we may as well take one for illustration. For example, for the first conditional log-likelihood, we can model it with EBM as:

$$p(\bar{z}_{i,\alpha}^u | \bar{z}_{i,\alpha}^v) = \frac{\exp(-E(\bar{z}_{i,\alpha}^u, \bar{z}_{i,\alpha}^v))}{\int \exp(-E(\bar{z}_{i,\alpha}^{u'}, \bar{z}_{i,\alpha}^v)) d\bar{z}_{i,\alpha}^{u'}} = \frac{\exp(f(\bar{z}_{i,\alpha}^u, \bar{z}_{i,\alpha}^v))}{A_{ij}}, \quad (6)$$

where $E(\cdot)$ is the energy function, A_{ij} is the intractable partition function, and $f(\cdot)$ is the negative energy. The energy function can be quite flexible, and for simplicity, we use the dot-product:

$$f(\bar{z}_{i,\alpha}^u, \bar{z}_{i,\alpha}^v) = \langle h(\mathbf{z}^u, \mathbf{d}_i, \alpha), h(\mathbf{z}^v, \mathbf{d}_i, \alpha) \rangle, \quad (7)$$

where $h(\cdot)$ is the editing function introduced in Eq. (2). Similarly for the other conditional log-likelihood term, and the objective becomes:

$$\mathcal{L}_{\text{GraphCG}} = \mathbb{E} \left[\log \frac{\exp(f(\bar{z}_{i,\alpha}^u, \bar{z}_{i,\alpha}^v))}{A_{ij}} + \log \frac{\exp(f(\bar{z}_{i,\alpha}^v, \bar{z}_{i,\alpha}^u))}{A_{ji}} \right]. \quad (8)$$

With Eq. (8), we are able to learn the semantically meaningful direction vectors. We name this unsupervised graph controllable generation framework as GraphCG. In specific, GraphCG utilizes EBM for estimation, which yields a wide family of solutions. Next we will introduce an intuitive solution.

4.3 GRAPHCG WITH NOISE CONTRASTIVE ESTIMATION

We solve Eq. (8) using the noise contrastive estimation (NCE) [16]. The high-level idea of NCE is to transform the density estimation problem into a binary classification problem that distinguishes if the data comes from the introduced noise distribution or from the true distribution. NCE has been widely explored for solving EBM [55], and we adopt it as GraphCG-NCE by optimizing:

$$\begin{aligned} \mathcal{L}_{\text{GraphCG-NCE}} = & - \left(\mathbb{E}_{p_n(\bar{z}_{j,\beta}^u | \bar{z}_{i,\alpha}^v)} [\log (1 - \sigma(f(\bar{z}_{j,\beta}^u, \bar{z}_{i,\alpha}^v)))] + \mathbb{E}_{p_{\text{data}}(\bar{z}_{i,\alpha}^u | \bar{z}_{i,\alpha}^v)} [\log \sigma(f(\bar{z}_{i,\alpha}^u, \bar{z}_{i,\alpha}^v))] \right. \\ & \left. + \mathbb{E}_{p_n(\bar{z}_{j,\beta}^v | \bar{z}_{i,\alpha}^u)} [\log (1 - \sigma(f(\bar{z}_{j,\beta}^v, \bar{z}_{i,\alpha}^u)))] + \mathbb{E}_{p_{\text{data}}(\bar{z}_{i,\alpha}^v | \bar{z}_{i,\alpha}^u)} [\log \sigma(f(\bar{z}_{i,\alpha}^v, \bar{z}_{i,\alpha}^u))] \right), \end{aligned} \quad (9)$$

where p_{data} is the data distribution and p_n is the noise distribution (derivations are in Appendix D). Recall that the latent pairs are given, and the noise distribution is on the semantic directions and step sizes. In specific, the step sizes ($\alpha \neq \beta$) are randomly sampled from $[-3, 3]$, and the latent direction indices ($i \neq j$) are randomly sampled from $\{0, 1, \dots, D-1\}$. The objective Eq. (9) is for one latent code pair, and we will take the expectation of it over all the pairs from the dataset. Besides, we would like to consider extra similarity and sparsity constraints as:

$$\mathcal{L}_{\text{sim}} = \mathbb{E}_{i,j} [\text{sim}(\mathbf{d}_i, \mathbf{d}_j)], \quad \mathcal{L}_{\text{sparsity}} = \mathbb{E}_i [\|\mathbf{d}_i\|], \quad (10)$$

where $\text{sim}(\cdot)$ is the similarity function between two latent directions, and we use the dot product. By minimizing these two regularization terms, we can make the learned semantic directions more diverse and sparse. Putting them together, the final objective function is:

$$\mathcal{L} = c_1 \cdot \mathbb{E}_{u,v} [\mathcal{L}_{\text{GraphCG-NCE}}] + c_2 \cdot \mathcal{L}_{\text{sim}} + c_3 \cdot \mathcal{L}_{\text{sparsity}}, \quad (11)$$

where c_1, c_2, c_3 are coefficients, and we treat them as three hyperparameters (check Appendix E). The above pipeline is illustrated in Fig. 1, and for the next we will discuss certain key modules.

Latent pairs, positive and negative views. We consider two options in designing the latent pairs. (1) *Perturbation (GraphCG-P)* is that for each latent variable $\mathbf{z} \in \mathcal{Z}$, we apply 2 perturbations (e.g., adding Gaussian noise) on \mathbf{z} to get 2 perturbed latent codes as \mathbf{z}^u and \mathbf{z}^v respectively. (2) *Random sampling (GraphCG-R)* is that we encode two randomly sampled data points from the empirical data distribution as \mathbf{z}^u and \mathbf{z}^v respectively. Perturbation is one of the widely-used strategies [28] for data augmentation, and random sampling has been widely used in the NCE [55] literature. Then we can define the positive and negative pairs in GraphCG-NCE, where the goal is to align the positives and contrast the negatives. As described in Eq. (9), the positive pairs are latent pairs moving with the same semantic direction and step size, while the negative pairs are the edited latent codes with different semantic directions and/or step sizes.

Semantic direction modeling. We first randomly draw a *basis vector* e_i , and then model the semantic direction d_i as $d_i = \text{MLP}(e_i)$, where $\text{MLP}(\cdot)$ is the multi-layer perceptron network.

Design of editing function. Given the semantic direction and two views, the next task is to design the editing function $h(\cdot)$ in Eq. (2). Since our proposed GraphCG is flexible, and the editing function determines the energy function Eq. (7), we consider both the linear and non-linear editing functions as:

$$\bar{z}_i = z + \alpha \cdot d_i, \quad \bar{z}_i = z + \alpha \cdot d_i + \text{MLP}(z \oplus d_i \oplus [\alpha]), \quad (12)$$

where \oplus is the concatenation of two vectors. Noticing that for the non-linear case, we are adding an extra term by mapping from the latent code, semantic direction, and step-size simultaneously. We expect that this could bring in more modeling expressiveness in the editing function. For more details, *e.g.*, the ablation study to check the effect on the design of the views and editing functions, please refer to [Appendices F and G](#), while more potential explorations are left for future work.

4.4 IMPLEMENTATIONS

During training, the goal of GraphCG is to learn semantically meaningful direction vectors together with an editing function in the latent space, as in Algorithm 1. Then we need to manually annotate the semantic directions with respect to the corresponding factors, using certain post-training evaluation metrics. Finally for the test phase, provided with the pretrained DGM on graph and a selected semantic direction (together with a step size), we can sample a molecule and use GraphCG for editing, as described in Eq. (2). The detailed algorithm is illustrated in Algorithm 2. Next, we highlight several key concepts in GraphCG and briefly discuss the differences to other related methods.

NCE and contrastive representation learning. GraphCG-NCE is applying EBM-NCE, which is essentially a contrastive learning method, and another dominant contrastive loss is the InfoNCE [42]. We summarize their relations as below. (1) Both contrastive methods are doing the same thing: align the positive pairs and contrast the negative pairs. (2) EBM-NCE [18, 35] has been found to outperform InfoNCE on the certain graph

applications like representation learning. (3) What we want to propose here is a flexible framework. Specifically, EBM provides a more general framework by designing the energy functions, and EBM-NCE is just one effective solution. Other promising directions include the denoising score matching or denoising diffusion model [56], while InfoNCE lacks such nice extensibility attribute.

GraphCG and self-supervised learning (SSL). GraphCG shares certain similarities with the self-supervised learning (SSL) method, however there are some inherent differences, as summarized below. (1) SSL aims at learning the data representation [by operating data augmentation on the data space, such as node addition and edge deletion](#). GraphCG aims at learning the semantically meaningful directions [by editing on the latent space \(the representation function is pretrained and fixed\)](#). (2) Based on the first point, SSL aims at using different data points as the negative samples. GraphCG, on the other hand, is using different directions and step-sizes as negatives. Namely, SSL is learning data representation in the inter-data level, and GraphCG is learning the semantic directions in the inter-direction level.

Algorithm 1 Training Phase of GraphCG

- 1: **Input:** Given a pretrained generative model, $f(\cdot)$ and $g(\cdot)$.
 - 2: **Output:** Learned direction vector d_i and function $h(\cdot)$.
 - 3: Select two data points, $z^u, z^v \in \mathcal{Z}$.
 - 4: **for** each step size α and each direction i **do**
 - 5: Set $\bar{z}_{i,\alpha}^u = h(z^u, d_i, \alpha)$.
 - 6: Set $\bar{z}_{i,\alpha}^v = h(z^v, d_i, \alpha)$.
 - 7: Assign positive to pair $(\bar{z}_{i,\alpha}^u, \bar{z}_{i,\alpha}^v)$.
 - 8: **for** step size $\beta \neq \alpha$ and direction $j \neq i$ **do**
 - 9: Set $\bar{z}_{j,\beta}^u = h(z^u, d_j, \beta)$.
 - 10: Set $\bar{z}_{j,\beta}^v = h(z^v, d_j, \beta)$.
 - 11: Assign negative to pair $(\bar{z}_{i,\alpha}^u, \bar{z}_{j,\beta}^v)$.
 - 12: Assign negative to pair $(\bar{z}_{j,\beta}^u, \bar{z}_{i,\alpha}^v)$.
 - 13: **end for**
 - 14: Do SGD w.r.t. GraphCG in Eq. (11).
 - 15: **end for**
-

Algorithm 2 Test Phase of GraphCG

- 1: **Input:** Given a pre-trained generative model ($f(\cdot)$ and $g(\cdot)$), a learned direction vector d .
 - 2: **Output:** A sequence of edited graphs.
 - 3: Sample a molecule with DGM or from a large molecule pool.
 - 4: Encode the molecule to get a latent code z .
 - 5: **for** step size $\alpha \in [-3, 3]$ **do**
 - 6: Do graph edit in the latent space to get $\bar{z}_{i,\alpha} = h(z, d, \alpha)$.
 - 7: Decode to the graph space with $\bar{x}' = g(\bar{z}_{i,\alpha})$.
 - 8: **end for**
 - 9: Output is thus a sequence of edited graphs, $\{\bar{x}'\}$.
-

Output sequence in the discrete space. Recall that during inference time (Algorithm 2), GraphCG takes a DGM and the learned semantic direction to output a sequence of edited graphs. Comparing to the vision domain, where certain models [51, 52] have proven their effectiveness in many tasks, the backbone models in the graph domain have limited discussions. This is challenging because the graph data is in a discrete and structured space, and the evaluation on such space is non-trivial. Meanwhile, GraphCG essentially provides another way to testify the quality of graph generation models. We would like to leave this for future exploration.

5 EXPERIMENTS

In this section, we show both the qualitative and quantitative results of GraphCG on two types of graph data: molecular graphs and point clouds. Due to the page limit, We put the experiment and implementation details in Appendix E. For reproducibility, the code will be public in the near future.

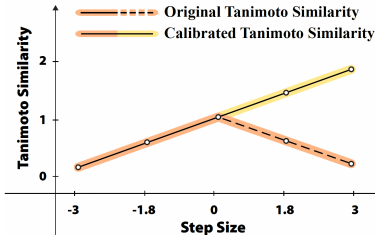
5.1 GRAPH DATA: MOLECULAR GRAPHS

Backbone DGMs. We consider two state-of-the-art DGMs for molecular graph generation. MoFlow [65] is a flow-based generative model on molecules which adopts an invertible mapping between the input molecular graphs and a latent prior. HierVAE [27] is a hierarchical VAE model which encodes and decodes molecule atoms and motifs in a hierarchical manner. Besides, the pretrained checkpoints are also provided, on ZINC250K [25] and ChEMBL [37] dataset respectively.

Editing sequences and anchor molecule. As discussed in Sec. 4, the output of the inference in GraphCG is a sequence of edited molecules with the i -th semantic direction, $\{\bar{x}'\}_i$. We first randomly generate a molecule using the backbone DGMs (without the editing operation), and we name such molecule as the anchor molecule, \bar{x}^* . Then we take 21 step sizes from -3 to 3, with interval 0.3, to obtain a sequence of 21 molecules following Eq. (2). Note that the edited molecule with step size 0 under the linear editing function is the same as the anchor molecule, *i.e.*, \bar{x}^* .

Change of structure factors and evaluation metrics. We are interested in the change of the graph structure (the steerable factors) along the output sequence edited with the i -th semantic direction. To evaluate the structure change, we apply the Tanimoto similarity between each output molecule and the anchor molecule. Besides, for the ease of evaluating the monotonicity, we apply a Tanimoto similarity transformation on the output molecules with positive step sizes by taking the deduction from 2. We call this calibrated Tanimoto similarity (CTS) sequence, marked as $\{s(\bar{x}')\}_i$. An illustration is shown in Fig. 2. Further, we propose a metric called Sequence Monotonic Ratio (SMR), $\phi_{\text{SMR}}(\gamma, \tau)_i$, which measures the monotonic ratio of M generated sequences edited with the i -th direction. It has two arguments: the diversity threshold γ constrains the minimum number of distinct molecules, and the tolerance threshold τ controls the non-monotonic tolerance ratio along each sequence.

Evaluating the diversity of semantic directions. SMR can evaluate the monotonic ratio of output sequences generated by one direction. To better illustrate that GraphCG is able to learn multiple directions with diverse semantic information, we also consider taking the average of top- K SMR to reveal that all the best K directions are semantically meaningful, as in Eq. (15).



$$\phi_{\text{SMR}}(\{s(\bar{x}')\}_i^m, \gamma, \tau) = \begin{cases} 1, & \text{len}(\text{set}(\{s(\bar{x}')\}_i^m)) \geq \gamma \\ & \wedge \text{monotonic}_{\tau}(\{s(\bar{x}')\}_i^m) \\ 0, & \text{otherwise} \end{cases}, \quad (13)$$

$$\phi_{\text{SMR}}(\gamma, \tau)_i = \frac{1}{M} \sum_{m=1}^M \phi_{\text{SMR}}(\{s(\bar{x}')\}_i^m, \gamma, \tau), \quad (14)$$

$$\text{top-K}(\gamma, \tau) = \frac{1}{K} \sum_{i \in \text{top-K directions}} (\phi_{\text{SMR}}(\gamma, \tau)_i). \quad (15)$$

Figure 2: This shows the sequence monotonic ratio (SMR) on calibrated Tanimoto similarity (CTS). Eqs. (13) and (14) are the SMR on each sequence and each direction respectively, where M is the number of generated sequences for the i -th direction and $\{s(\bar{x}')\}_i^m$ is the CTS of the m -th generated sequence with the i -th direction. Eq. (15) is the average of top- K SMR on D directions. More details are in Appendix F.

Table 2: This table lists the sequence monotonic ratio (SMR, %) on calibrated Tanimoto similarity (CTS) with respect to the top-1 and top-3 directions. The best performances are marked in **bold**.

Model	Dataset	diversity γ tolerance τ	Tanimoto top-1				Tanimoto top-3			
			3		4		3		4	
			0	0.2	0	0.2	0	0.2	0	0.2
MoFlow	ZINC250k	Random	23.0	25.0	12.0	15.0	22.0	24.0	11.0	13.7
		Variance	24.0	28.0	12.0	16.0	20.0	25.0	10.0	15.0
		SeFa [51]	4.0	4.0	0.0	0.0	3.3	3.3	0.0	0.0
		DisCo [45]	7.0	14.0	2.0	8.0	5.3	11.7	2.0	7.7
		GraphCG-P	32.0	34.0	16.0	18.0	29.0	31.0	13.7	16.3
		GraphCG-R	25.0	26.0	11.0	14.0	22.0	24.3	10.3	13.3
HierVAE	ChEMBL	Random	14.0	45.0	14.0	43.0	10.0	42.3	9.3	41.7
		Variance	23.0	59.0	19.0	55.0	18.3	52.7	15.7	50.3
		SeFa [51]	4.0	41.0	4.0	41.0	2.3	36.0	2.3	36.0
		GraphCG-P	40.0	73.0	32.0	65.0	36.0	64.3	26.3	57.7
		GraphCG-R	42.0	67.0	30.0	55.0	38.0	62.3	28.7	53.7

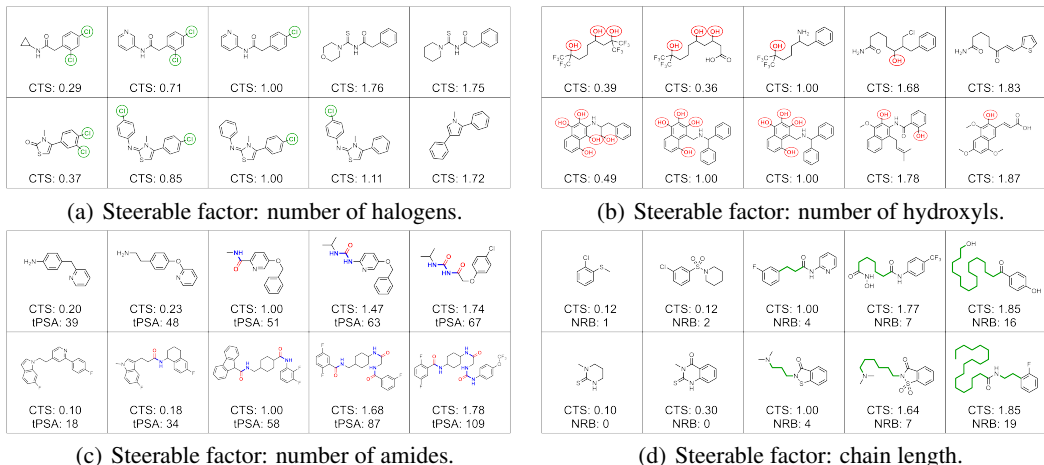


Figure 3: GraphCG for molecular graph editing. We visualize the output molecules and CTS on four directions with two sequences each, where each sequence consists of five steps. The center point is the anchor molecule, and the other four points correspond to step size with -3, -1.8, 1.8, and 3 respectively. Figs. 3(a) to 3(c) show how functional groups in the molecules can be viewed as the steerable factors as they change along the sequence, such as halogen atoms, hydroxyl groups and amides. Fig. 3(d) illustrates the effect on the steerable factor on the length of flexible chains in the molecules. Notably, certain properties change with molecular structures accordingly, like topological polar surface area (tPSA) and number of rotatable bonds (NRB).

Baselines. For baselines, we consider four unsupervised editing methods. (1) The first is Random. It randomly samples a normalized random vector in the latent space as the semantic direction. (2) The second one is Variance. We analyze the variance on each dimension of the latent space, and select the highest one with one-hot encoding as the semantic direction. (3) The third one is SeFa [51]. It first decomposes the latent space into lower dimensions using PCA, and then takes the most principle components (eigenvectors) as the semantic-rich direction vectors. (4) The last one is DisCo [45]. It maps each latent code back to the data space, followed with an encoder for contrastive learning. Namely, it requires the backbone DGMs to be end-to-end and is infeasible for HierVAE.

Quantitative results. We take $D = 10$ to train GraphCG, and the optimal results on 100 sampled sequences are reported in Table 2. We can observe that GraphCG can show consistently better structure change with both top-1 and top-3 directions. This can empirically prove the effectiveness of our proposed GraphCG. More comprehensive results are in Appendix F.

Analysis on steerable factors in molecules: functional group change. For visualization, we sample 8 molecular graph sequences along 4 selected directions in Fig. 3, and the backbone DGM is HierVAE pretrained on ChEMBL. The CTS holds good monotonic trend, and each direction shows certain unique changes in the molecular structures, *i.e.*, the steerable factors in molecules. Some structural changes are reflected in molecular properties. We expand all the details below. In Fig. 3(a) and Fig. 3(b), the number of halogen atoms and hydroxyl groups (in alcohols and phenols) in the

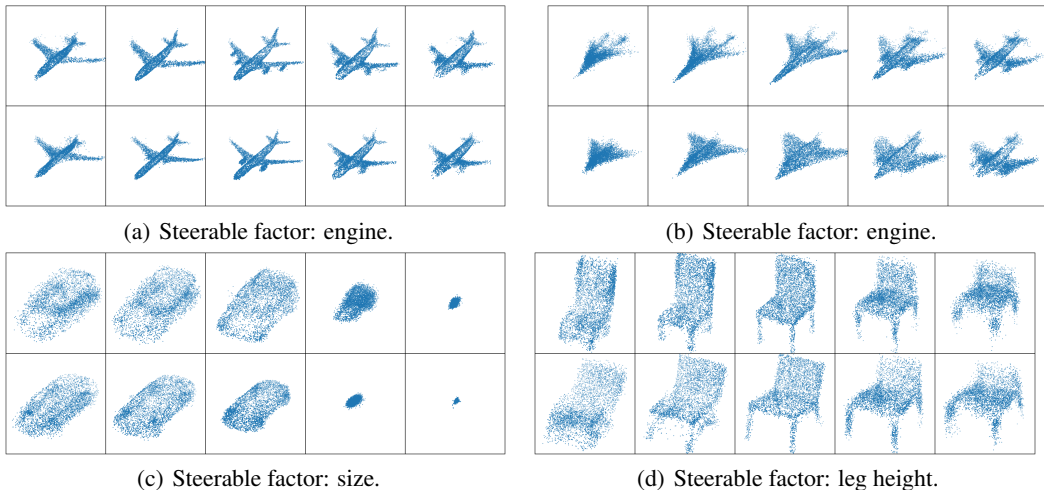


Figure 4: GraphCG for point clouds editing. We show four editing sequences, where each sequence consists of five point clouds, and the center one is the anchor point clouds, *i.e.*, with step size 0. The other four point clouds correspond to step size with -3, -1.8, 1.8, and 3, respectively. Fig. 4(a) and Fig. 4(b) refer to the same semantic direction, and they are showing how to steer the factor engine: the number of engines will be decreased and increased with the negative (left) and positive (right) step size respectively. Similarly, Figs. 4(c) and 4(d) illustrate the effect on the steerable factors on the car size and the chair leg height.

molecules decrease from left to right, respectively. In Fig. 3(c), the number of amides in the molecules increases along the path. Because amides are polar functional groups, the topological polar surface area (tPSA) of the molecules also increase accordingly, which is a key molecular property for the prediction of drug transport properties, *e.g.*, permeability [12]. In Fig. 3(d), the flexible chain length, marked by the number of ethylene (CH_2CH_2) units, increases from left to right. Since the number of rotatable bonds (NRB) measures the molecular flexibility, it also increases accordingly [57].

5.2 GRAPH DATA: POINT CLOUDS

Backbone DGMs. We consider one of the latest DGMs on point clouds, PointFlow [62]. It is using the normalizing flow model for estimating the 3D point clouds distribution. Then we consider PointFlow pretrained on three datasets in ShapeNet [3]: Airplane, Car, and Chair. All point clouds are obtained by sampling points uniformly from the mesh surface.

Analysis on steerable factors in point clouds: shape change. To train GraphCG, we take $D = 10$ directions, and we sample 8 point cloud sequences along 3 directions for visualization in Fig. 4. More results are in Appendix G. It is observed that GraphCG can steer the shape of the point clouds, *e.g.*, the size of cars and the height of chair legs. We also find it interesting that GraphCG can steer more finger-trained factors, like modifying the number jet engines of airplanes in Figs. 4(a) and 4(b).

6 CONCLUSION AND DISCUSSION

In this work, we are interested in unsupervised graph editing. It is a well-motivated task for various real-world applications, and we discuss two mainstream data types: molecular graphs and point clouds. We start with exploring the latent space of mainstream deep generative models and propose GraphCG, a model-agnostic and task-agnostic unsupervised method for graph data editing. The key component of GraphCG is EBM, and we take the GraphCG-NCE as the solution for now. For the future work, we may as well extend it to more advanced solutions like denoising diffusion model [21].

One limitation of GraphCG (as well as the solutions to general unsupervised data editing) [17, 45, 51] is that we may need some post-training human selection (as shown in Algorithms 1 and 2) to select the most promising semantic vectors to steer factors. Another issue is the lack of open-sourced evaluation metrics. This requires both a deep understanding of the representation space of deep generative models and domain knowledge of the problem. For instance, activity cliff is a challenging task [22] for editing, while current measures fail to capture it. To set up constructive evaluation metrics can help augment our understandings from both the domain and technique perspectives. This is beyond the scope of our work, yet would be interesting to explore as a future direction.

ETHICS STATEMENT

We authors acknowledge that we have read and commit to adhering to the ICLR Code of Ethics.

REPRODUCIBILITY STATEMENT

To ensure the reproducibility of the empirical results, we provide the implementation details (hyperparameters, dataset specifications, pretrained checkpoints, etc.) in Sec. 5 and appendices C and E to G, and the source code will be released in the future. Besides, the complete derivations of equations and clear explanations are given in Sec. 4 and appendix D.

Specifically, we provide the details for reproducing the results:

- In Table 5, GraphCG-P with Eq. (26) and GraphCG-R with Eq. (24) are reported in Table 2.
- In Table 6, GraphCG-P with Eq. (25) and GraphCG-R with Eq. (25) are reported in Table 2.

For the visualization in Fig. 3, we take the GraphCG-P with Eq. (25), and the backbone generative model is HierVAE pretrained on ChEMBL. Further, we provide an anonymous link [here](#). In these CSV files:

- Direction 0 is the halogen fragment (data 4, 71).
- Direction 5 is the amide fragment (data 95, 61).
- Direction 6 is the chain length (data 57, 14).
- Direction 7 is the alcohol and phenol fragments (data 10, 8).

For the visualization in Fig. 4, we take the GraphCG-R with Eq. (25) on PointFlow, for all three datasets.

REFERENCES

- [1] Michael Arbel, Liang Zhou, and Arthur Gretton. Generalized energy based models. *arXiv:2003.05033*, 2020. 17
- [2] Hans-Joachim Böhm, Alexander Flohr, and Martin Stahl. Scaffold hopping. *Drug discovery today: Technologies*, 1(3):217–224, 2004. 1
- [3] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, et al. Shapenet: An information-rich 3d model repository. *arXiv:1512.03012*, 2015. 9
- [4] Tong Che, Ruixiang Zhang, Jascha Sohl-Dickstein, Hugo Larochelle, Liam Paull, et al. Your gan is secretly an energy-based model and you should use discriminator driven latent sampling. *arXiv:2003.06060*, 2020. 3
- [5] Hongming Chen, Ola Engkvist, Yinhai Wang, Marcus Olivecrona, and Thomas Blaschke. The rise of deep learning in drug discovery. *Drug discovery today*, 23(6):1241–1250, 2018. 2
- [6] Tian Qi Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *NeurIPS*, 2018. 4, 15
- [7] Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal neighbourhood aggregation for graph nets. *arXiv:2004.05718*, 2020. 14
- [8] Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. *arXiv:1410.8516*, 2014. 3
- [9] Jürgen Drews. Drug discovery: a historical perspective. *Science*, 287(5460):1960–1964, 2000. 1, 14
- [10] David Duvenaud, Dougal Maclaurin, Jorge Aguilera-Iparraguirre, Rafael Gómez-Bombarelli, Timothy Hirzel, et al. Convolutional networks on graphs for learning molecular fingerprints. *arXiv:1509.09292*, 2015. 1, 14
- [11] Cian Eastwood and Christopher KI Williams. A framework for the quantitative evaluation of disentangled representations. In *ICLR*, 2018. 1, 4, 15
- [12] Peter Ertl, Bernhard Rohde, and Paul Selzer. Fast calculation of molecular polar surface area as a sum of fragment-based contributions and its application to the prediction of drug transport properties. *Journal of medicinal chemistry*, 43(20):3714–3717, 2000. 9, 21

- [13] Peter Ertl, Eva Altmann, and Jeffrey M McKenna. The most common functional groups in bioactive molecules and how their popularity has evolved over time. *Journal of medicinal chemistry*, 63(15): 8408–8418, 2020. 1
- [14] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *ICML*, 2017. 1, 14
- [15] Laurent Gomez. Decision making in medicinal chemistry: The power of our intuition. *ACS Medicinal Chemistry Letters*, 9(10):956–958, 2018. 1
- [16] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *AISTATS*, 2010. 3, 5
- [17] Erik Härkönen, Aaron Hertzman, Jaakko Lehtinen, and Sylvain Paris. Ganspace: Discovering interpretable gan controls. In *NeurIPS*, 2020. 2, 3, 9, 14
- [18] Kaveh Hassani and Amir Hosein Khasahmadi. Contrastive multi-view representation learning on graphs. In *ICML*, 2020. 6
- [19] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, et al. beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2017. 1, 3, 4, 15
- [20] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002. 3
- [21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 9
- [22] Ye Hu and Bajorath Jurgen. Extending the activity cliff concept: structural categorization of activity cliffs and systematic identification of different types of cliffs in the chembl database. *Journal of chemical information and modeling*, 52(7):1806–1811, 2012. 9
- [23] Ye Hu, Dagmar Stumpfe, and Jurgen Bajorath. Recent advances in scaffold hopping: miniperspective. *Journal of medicinal chemistry*, 60(4):1238–1246, 2017. 1
- [24] Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005. 3
- [25] John J Irwin and Brian K Shoichet. Zinc: a free database of commercially available compounds for virtual screening. *Journal of chemical information and modeling*, 45(1):177–182, 2005. 7
- [26] Ali Jahanian, Lucy Chai, and Phillip Isola. On the "steerability" of generative adversarial networks. In *ICLR*, 2019. 2, 14
- [27] Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Hierarchical generation of molecular graphs using structural motifs. In *ICML*, 2020. 2, 7, 14
- [28] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 3, 5, 14
- [29] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, et al. Analyzing and improving the image quality of stylegan. In *CVPR*, 2020. 3
- [30] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *ICML*, 2018. 4, 15
- [31] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv:1312.6114*, 2013. 3
- [32] Abhishek Kumar, Prasanna Sattigeri, and Avinash Balakrishnan. Variational inference of disentangled latent concepts from unlabeled observations. In *ICLR*, 2018. 4, 15
- [33] Greg Landrum et al. RDKit: A software suite for cheminformatics, computational chemistry, and predictive modeling, 2013. 4
- [34] Shengchao Liu, Mehmet Furkan Demirel, and Yingyu Liang. N-gram graph: Simple unsupervised representation for graphs, with applications to molecules. *arXiv:1806.09206*, 2018. 14
- [35] Shengchao Liu, Hanchen Wang, Weiyang Liu, Joan Lasenby, Hongyu Guo, et al. Pre-training molecular graph representation with 3d geometry. In *ICLR*, 2022. 4, 6, 16, 17

- [36] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, et al. Challenging common assumptions in the unsupervised learning of disentangled representations. In *ICML*, 2019. 1, 2, 3, 15
- [37] David Mendez, Anna Gaulton, A Patrícia Bento, Jon Chambers, Marleen De Veij, et al. ChEMBL: towards direct deposition of bioassay data. *Nucleic acids research*, 47(D1):D930–D940, 2019. 7
- [38] Zlatko Mihalić and Nenad Trinajstić. A graph-theoretical approach to structure-property relationships. *Journal of Chemical Education*, 69(9):701, 1992. 1
- [39] Andriy Mnih and Yee Whye Teh. A fast and simple algorithm for training neural probabilistic language models. *arXiv preprint arXiv:1206.6426*, 2012. 17
- [40] Weili Nie, Arash Vahdat, and Anima Anandkumar. Controllable and compositional generation with latent-space energy-based models. In *NeurIPS*, 2021. 2, 14
- [41] Marcus Olivecrona, Thomas Blaschke, Ola Engkvist, and Hongming Chen. Molecular de-novo design through deep reinforcement learning. *Journal of cheminformatics*, 9(1):1–14, 2017. 14
- [42] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv:1807.03748*, 2018. 6
- [43] Yael Pritch, Eitam Kav-Venaki, and Shmuel Peleg. Shift-map image editing. In *ICCV*, 2009. 1
- [44] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, 2017. 14
- [45] Xuanchi Ren, Tao Yang, Yuwang Wang, and Wenjun Zeng. Do generative models know disentanglement? contrastive learning is all you need. *arXiv:2102.10543*, 2021. 2, 8, 9, 14
- [46] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *ICML*, 2015. 3
- [47] Karl Ridgeway and Michael C Mozer. Learning deep disentangled embeddings with the f-statistic loss. In *NeurIPS*, 2018. 1, 4, 15
- [48] Radu Bogdan Rusu and Steve Cousins. 3D is here: Point Cloud Library (PCL). In *ICRA*, 2011. 4, 15
- [49] Radu Bogdan Rusu, Gary Bradski, Romain Thibaux, and John Hsu. Fast 3d recognition and pose using the viewpoint feature histogram. In *IROS*, 2010. 4, 15
- [50] Paul G. Seybold, Michael May, and Ujjvala A. Bagal. Molecular structure: Property relationships. *Journal of Chemical Education*, 64(7):575, 1987. 4
- [51] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in gans. In *CVPR*, 2021. 2, 3, 7, 8, 9, 14
- [52] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. Interfacegan: Interpreting the disentangled face representation learned by gans. *IEEE TPAMI*, 2020. 2, 3, 7, 14
- [53] Weijing Shi and Raj Rajkumar. Point-gnn: Graph neural network for 3d object detection in a point cloud. In *CVPR*, 2020. 1, 14
- [54] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *arXiv:1907.05600*, 2019. 3
- [55] Yang Song and Diederik P Kingma. How to train your energy-based models. *arXiv:2101.03288*, 2021. 5, 17
- [56] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, et al. Score-based generative modeling through stochastic differential equations. *arXiv:2011.13456*, 2020. 3, 6
- [57] Daniel F Veber, Stephen R Johnson, Hung-Yuan Cheng, Brian R Smith, Keith W Ward, and Kenneth D Kopple. Molecular properties that influence the oral bioavailability of drug candidates. *Journal of medicinal chemistry*, 45(12):2615–2623, 2002. 2, 9, 21
- [58] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, et al. Dynamic graph cnn for learning on point clouds. *ACM ToG*, 2020. 1
- [59] Zichao Wang, Weili Nie, Zhuoran Qiao, Chaowei Xiao, Richard Baraniuk, and Anima Anandkumar. Retrieval-based controllable molecule generation. *arXiv preprint arXiv:2208.11126*, 2022. 14

- [60] Zongze Wu, Dani Lischinski, and Eli Shechtman. Stylespace analysis: Disentangled controls for stylegan image generation. In *CVPR*, 2021. 1
- [61] Weihao Xia, Yulun Zhang, Yujiu Yang, Jing-Hao Xue, Bolei Zhou, et al. Gan inversion: A survey. *arXiv:2101.05278*, 2021. 2
- [62] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, et al. Pointflow: 3d point cloud generation with continuous normalizing flows. In *ICCV*, 2019. 9, 14
- [63] Kevin Yang, Kyle Swanson, Wengong Jin, Connor Coley, Philipp Eiden, et al. Analyzing learned molecular representations for property prediction. *Journal of chemical information and modeling*, 59(8):3370–3388, 2019. 14
- [64] Jiaxuan You, Bowen Liu, Zhitao Ying, Vijay Pande, and Jure Leskovec. Graph convolutional policy network for goal-directed molecular graph generation. *Advances in neural information processing systems*, 31, 2018. 2, 14
- [65] Chengxi Zang and Fei Wang. Moflow: an invertible flow model for generating molecular graphs. In *KDD*, 2020. 7, 14

A GRAPH DATA

A.1 MOLECULES

Molecules can be naturally represented as the 2D molecular graphs, where atoms and bonds are nodes and edges in the graph, respectively. For the recent years, graph representation learning has been extensively explored on the molecular graph [7, 10, 14, 34, 63]. Based on the graph representation of molecules, a variety of work have been done for molecule generation. The state-of-the-art ones include MoFlow [65] and HierVAE [27].

A.2 POINT CLOUDS

A point cloud is represented as a set of points, where each point P_i is described by a vector of 3D Euclidean coordinates possibly with extra channels (e.g., colors, surface normals and returned laser intensity). In 2017, Qi et. al [44] designed a deep learning framework called PointNet that directly consumes unordered point sets as inputs and can be used for various tasks such as classification and segmentation. For the generative models on point clouds, we consider one of the latest work, PointFlow [62].

B RELATED WORK

Image editing and image controllable generation Many existing works on controllable generation with DGMs mainly focus on the image data. With the assumption that the learned latent space already include rich semantic information [17, 26, 28, 52], the question then becomes how to identify semantic-rich directions from the latent space of DGMs. Depending on whether or using supervised signals to discover the semantic directions, existing works can be divided into two settings: *supervised* and *unsupervised*.

The *supervised* setting relies on the supervised signals to learn the pre-defined semantic-rich directions. For instance, InterfaceGAN [52] identifies the semantic directions in the latent space via linear models that recognize semantic boundary among data. LACE [40] uses energy-based models to learn the joint distribution of data and properties (*i.e.*, semantics) and formulate the sampling process as to solve an ordinary differential equation.

As supervised signals usually require domain knowledge and laborious annotations, latest works are more focused on the *unsupervised* setting, either model-specific or data-specific. Specifically, GANSpace [17] applies PCA on the intermediate layers of GANs (instead of latent space) for learning the semantic-rich directions. SeFa [51] exploits the pretrained GANs and extracts the semantic-rich directions by using PCA on the backbone layers. Nevertheless, as both methods are specifically designed for StyleGAN [28], it is nontrivial to generalize them to other DGMs. A more recent work DisCo [45] employs a different pipeline: it trains a new encoder after reconstruction, and maps the generated images to another latent space for editing. However, training a new encoder introduces extra complexities.

Graph editing and graph controllable generation This is an emerging field with many downstream applications [9, 53]. However, existing works are mainly focusing on the supervised setting. For example, conventional approaches, such as genetic algorithms (GAs), edit the molecule graphs in the data space via hand-crafted heuristics with the guidance of molecular property predictors. More recent learning based methods perform the latent direction discovery, either by training a classifier on the latent space of DGMs on the graph data [27] or by learning to retrieve exemplar samples from a retrieval database for guidance [59]. Other works fine-tune a pre-trained graph DGM using the supervised property annotations as rewards, resulting a controllable DGM specifically for the considered task [41, 64]. To the best of our knowledge, our work is the first to explore the unsupervised graph editing in the unsupervised manner. Besides, different from many previous approaches that may only work for molecule graphs or point cloud graphs, our method is generic and thus can be applied to various graph modalities.

C ANALYSIS EXPERIMENTS ON DISENTANGLEMENT

In Sec. 3, we conduct three analysis experiments to conclude that the representation space is not perfectly disentangled in such the setting. In this section, we provide more details and complementary information of these experiments.

C.1 STEERABLE FACTORS

As mentioned in Sec. 3, we consider the measuring the disentanglement with respect to three types of structured data: molecular graphs and point clouds. Recall that we need to define steerable factors first, so as to measure the disentanglement.

Table 3: The six mainstream disentanglement metrics on five DGMs and three data types. All measures range from 0 to 1, and higher scores mean more disentangled representation. MoFlow and HierVAE are for molecular graphs, PointFlow is for point clouds.

Data Type	Model	Dataset	BetaVAE \uparrow	FactorVAE \uparrow	MIG \uparrow	DCI \uparrow	Modularity \uparrow	SAP \uparrow
Molecular Graphs	MolFlow	QM9	0.251	0.165	0.038	0.727	0.599	0.017
		ZINC250k	0.264	0.175	0.030	0.958	0.620	0.009
	HierVAE	QM9	0.165	0.135	0.044	0.241	0.626	0.076
		ChEMBL	0.159	0.130	0.023	0.113	0.604	0.026
Point Clouds	PointFlow	Airplane	0.022	0.025	0.029	0.160	0.745	0.022
		Chair	0.019	0.014	0.032	0.149	0.721	0.019
		Car	0.019	0.023	0.021	0.120	0.713	0.021

Molecular Graph For molecular graphs, we treat the important substructures (a.k.a., motifs or fragments) as factors, and they are extracted using RDKit. We consider the existence of the following 9 motifs as the binary factor labels:

- aliphatic hydroxyl groups.
- aliphatic hydroxyl groups excluding tert-OH.
- aromatic nitrogens.
- aromatic amines.
- Tertiary amines.
- Secondary amines.
- amides.
- ether oxygens (including phenoxy).
- nitriles.

Point Clouds For point clouds, we adopt the viewpoint feature histogram (VFH) [49] implemented in PCL [48]. There are 308 bins in total, where each bin is a histogram of the angles that viewpoint direction makes with each normal. VFH has been widely used as point cloud descriptors, and here we binarize it into factors with:

- We collect the shared non-zero bins from all three datasets (Airplane, Car, and Chair), and ignore the bins where the values distribution are highly skewed. This can give us 75 bins.
- Then for each of these selected bins, we use the median value as the threshold to generate the binary factor labels.

C.2 DISENTANGLEMENT MEASURES

We follow the [36] on considering the following six disentanglement measures:

- β -VAE [19] evaluates the prediction accuracy of a linear classifier for the index of a fixed factor of variation.
- FactorVAE [30] addresses the limitations (i.e. corner case) of β -VAE via introducing a majority voting classifier on a different feature vector.
- MIG [6] measures the normalized difference between the highest and second highest mutual information between latent dimensions and each steerable factor.
- DCI [11] disentanglement score measures the average difference from one of the entropy of probability that a latent dimension is useful for predicting a steerable factor (computed by the relative importance score).
- Modularity [47] measures whether each latent dimension is dependent on at most one single steerable factor. It computes the average normalized squared difference over the highest and second highest mutual information between each steerable factor and each latent dimension.
- SAP [32] calculates the R^2 score of the linear models trained to predict each steerable factor from each latent dimension.

Recall that all six disentanglement measures range from 0 to 1, and higher value means the corresponding space is more disentangled. The results can be found in Table 3. We can conclude that all the values are indeed low on all datasets and models, revealing that DGMs are not well-disentangled in general.

D MUTUAL INFORMATION

We added the comprehensive derivations below. In this section, we will briefly introduce mutual information (MI), and also a lower bound for maximizing MI. This has been previously proposed in [35] for self-supervised learning. First for notation, without loss of generality, we use X and Y as two views.

D.1 FORMULATION

The standard formulation for mutual information (MI) is

$$I(X; Y) = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} \left[\log \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} \right]. \quad (16)$$

Mutual information (MI) between random variables measures the corresponding non-linear dependence. As can be seen in the first equation in Eq. (16), the larger the divergence between the joint ($p(\mathbf{x}, \mathbf{y})$) and the product of the marginals $p(\mathbf{x})p(\mathbf{y})$, the stronger the dependence between X and Y .

D.2 A LOWER BOUND TO MI

First we can get a lower bound of MI. Assuming that there exist (possibly negative) constants a and b such that $a \leq H(X)$ and $b \leq H(Y)$, i.e., the lower bounds to the (differential) entropies, then we have:

$$\begin{aligned} I(X; Y) &= \frac{1}{2} (H(X) + H(Y) - H(Y|X) - H(X|Y)) \\ &\geq \frac{1}{2} (a + b - H(Y|X) - H(X|Y)) \\ &\geq \frac{1}{2} (a + b) + \mathcal{L}_{\text{MI}}. \end{aligned} \quad (17)$$

Thus, we transform the MI maximization problem into maximizing the following objective:

$$\mathcal{L}_{\text{MI}} = \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log p(\mathbf{x}|\mathbf{y})] + \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [\log p(\mathbf{y}|\mathbf{x})]. \quad (18)$$

Empirically, we use energy-based models to model the distributions. The existence of a and b can be understood as the requirements that the two distributions ($p_{\mathbf{x}}, p_{\mathbf{y}}$) are not collapsed. For the next, we will try to model the two conditional data distributions p with model distributions p_{θ} .

D.3 DERIVATION OF CONDITIONAL EBM WITH NCE

WLOG, let's consider modeling the $p_{\theta}(\mathbf{x}|\mathbf{y})$ first, and by EBM it is as follows:

$$p_{\theta}(\mathbf{x}|\mathbf{y}) = \frac{\exp(-E(\mathbf{x}|\mathbf{y}))}{\int \exp(-E(\tilde{\mathbf{x}}|\mathbf{y})) d\tilde{\mathbf{x}}} = \frac{\exp(f_{\mathbf{x}}(\mathbf{x}, \mathbf{y}))}{\int \exp(f_{\mathbf{x}}(\tilde{\mathbf{x}}|\mathbf{y})) d\tilde{\mathbf{x}}} = \frac{\exp(f_{\mathbf{x}}(\mathbf{x}, \mathbf{y}))}{A_{\mathbf{x}|\mathbf{y}}}. \quad (19)$$

Then we solve this using NCE. NCE handles the intractability issue by transforming it as a binary classification task. We take the partition function $A_{\mathbf{x}|\mathbf{y}}$ as a parameter, and introduce a noise distribution p_n . Based on this, we introduce a mixture model: $z = 0$ if the conditional $\mathbf{x}|\mathbf{y}$ is from $p_n(\mathbf{x}|\mathbf{y})$, and $z = 1$ if $\mathbf{x}|\mathbf{y}$ is from $p_{\text{data}}(\mathbf{x}|\mathbf{y})$. So the joint distribution is:

$$p_{n, \text{data}}(\mathbf{x}|\mathbf{y}) = p(z = 1)p_{\text{data}}(\mathbf{x}|\mathbf{y}) + p(z = 0)p_n(\mathbf{x}|\mathbf{y})$$

The posterior of $p(z = 0|\mathbf{x}, \mathbf{y})$ is

$$p_{n, \text{data}}(z = 0|\mathbf{x}, \mathbf{y}) = \frac{p(z = 0)p_n(\mathbf{x}|\mathbf{y})}{p(z = 0)p_n(\mathbf{x}|\mathbf{y}) + p(z = 1)p_{\text{data}}(\mathbf{x}|\mathbf{y})} = \frac{\nu \cdot p_n(\mathbf{x}|\mathbf{y})}{\nu \cdot p_n(\mathbf{x}|\mathbf{y}) + p_{\text{data}}(\mathbf{x}|\mathbf{y})},$$

where $\nu = \frac{p(z=0)}{p(z=1)}$.

Similarly, we can have the joint distribution under EBM framework as:

$$p_{n, \theta}(\mathbf{x}) = p(z = 0)p_n(\mathbf{x}|\mathbf{y}) + p(z = 1)p_{\theta}(\mathbf{x}|\mathbf{y})$$

And the corresponding posterior is:

$$p_{n, \theta}(z = 0|\mathbf{x}, \mathbf{y}) = \frac{p(z = 0)p_n(\mathbf{x}|\mathbf{y})}{p(z = 0)p_n(\mathbf{x}|\mathbf{y}) + p(z = 1)p_{\theta}(\mathbf{x}|\mathbf{y})} = \frac{\nu \cdot p_n(\mathbf{x}|\mathbf{y})}{\nu \cdot p_n(\mathbf{x}|\mathbf{y}) + p_{\theta}(\mathbf{x}|\mathbf{y})}$$

We indirectly match $p_\theta(\mathbf{x}|\mathbf{y})$ to $p_{\text{data}}(\mathbf{x}|\mathbf{y})$ by fitting $p_{n,\theta}(\mathbf{z}|\mathbf{x}, \mathbf{y})$ to $p_{n,\text{data}}(\mathbf{z}|\mathbf{x}, \mathbf{y})$ by minimizing their KL-divergence:

$$\begin{aligned}
& \min_{\theta} D_{\text{KL}}(p_{n,\text{data}}(\mathbf{z}|\mathbf{x}, \mathbf{y}) || p_{n,\theta}(\mathbf{z}|\mathbf{x}, \mathbf{y})) \\
&= \mathbb{E}_{p_{n,\text{data}}(\mathbf{x}, \mathbf{z}|\mathbf{y})} [\log p_{n,\theta}(\mathbf{z}|\mathbf{x}, \mathbf{y})] \\
&= \int \sum_{\mathbf{z}} p_{n,\text{data}}(\mathbf{x}, \mathbf{z}|\mathbf{y}) \cdot \log p_{n,\theta}(\mathbf{z}|\mathbf{x}, \mathbf{y}) d\mathbf{x} \\
&= \int \left\{ p(\mathbf{z} = 0) p_{n,\text{data}}(\mathbf{x}|\mathbf{y}, \mathbf{z} = 0) \log p_{n,\theta}(\mathbf{z} = 0|\mathbf{x}, \mathbf{y}) \right. \\
&\quad \left. + p(\mathbf{z} = 1) p_{n,\text{data}}(\mathbf{x}|\mathbf{y}, \mathbf{z} = 1) \log p_{n,\theta}(\mathbf{z} = 1|\mathbf{x}, \mathbf{y}) \right\} d\mathbf{x} \\
&= \nu \cdot \mathbb{E}_{p_n(\mathbf{x}|\mathbf{y})} \left[\log p_{n,\theta}(\mathbf{z} = 0|\mathbf{x}, \mathbf{y}) \right] + \mathbb{E}_{p_{\text{data}}(\mathbf{x}|\mathbf{y})} \left[\log p_{n,\theta}(\mathbf{z} = 1|\mathbf{x}, \mathbf{y}) \right] \\
&= \nu \cdot \mathbb{E}_{p_n(\mathbf{x}|\mathbf{y})} \left[\log \frac{\nu \cdot p_n(\mathbf{x}|\mathbf{y})}{\nu \cdot p_n(\mathbf{x}|\mathbf{y}) + p_\theta(\mathbf{x}|\mathbf{y})} \right] + \mathbb{E}_{p_{\text{data}}(\mathbf{x}|\mathbf{y})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{y})}{\nu \cdot p_n(\mathbf{x}|\mathbf{y}) + p_\theta(\mathbf{x}|\mathbf{y})} \right].
\end{aligned} \tag{20}$$

This optimal distribution is an estimation to the actual distribution (or data distribution), *i.e.*, $p_\theta(\mathbf{x}|\mathbf{y}) \approx p_{\text{data}}(\mathbf{x}|\mathbf{y})$. We can follow the similar steps for $p_\theta(\mathbf{y}|\mathbf{x}) \approx p_{\text{data}}(\mathbf{y}|\mathbf{x})$. Thus following Eq. (20), the objective function is to maximize

$$\nu \cdot \mathbb{E}_{p_{\text{data}}(\mathbf{y})} \mathbb{E}_{p_n(\mathbf{x}|\mathbf{y})} \left[\log \frac{\nu \cdot p_n(\mathbf{x}|\mathbf{y})}{\nu \cdot p_n(\mathbf{x}|\mathbf{y}) + p_\theta(\mathbf{x}|\mathbf{y})} \right] + \mathbb{E}_{p_{\text{data}}(\mathbf{y})} \mathbb{E}_{p_{\text{data}}(\mathbf{x}|\mathbf{y})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{y})}{\nu \cdot p_n(\mathbf{x}|\mathbf{y}) + p_\theta(\mathbf{x}|\mathbf{y})} \right]. \tag{21}$$

The we will adopt three strategies to approximate Eq. (21):

1. **Self-normalization.** When the EBM is very expressive, *i.e.*, using deep neural network for modeling, we can assume it is able to approximate the normalized density directly [39, 55]. In other words, we can set the partition function $A = 1$. This is a self-normalized EBM-NCE, with normalizing constant close to 1, *i.e.*, $p(\mathbf{x}) = \exp(-E(\mathbf{x})) = \exp(f(\mathbf{x}))$.
2. **Exponential tilting term.** Exponential tilting term [1] is another useful trick. It models the distribution as $\tilde{p}_\theta(\mathbf{x}) = q(\mathbf{x}) \exp(-E_\theta(\mathbf{x}))$, where $q(\mathbf{x})$ is the reference distribution. If we use the same reference distribution as the noise distribution, the tilted probability is $\tilde{p}_\theta(\mathbf{x}) = p_n(\mathbf{x}) \exp(-E_\theta(\mathbf{x}))$.
3. **Sampling.** For many cases [35], we only need to sample 1 negative points for each data, *i.e.*, $\nu = 1$.

Following these three disciplines, the objective function to optimize $p_\theta(\mathbf{x}|\mathbf{y})$ becomes

$$\begin{aligned}
& \mathbb{E}_{p_n(\mathbf{x}|\mathbf{y})} \left[\log \frac{p_n(\mathbf{x}|\mathbf{y})}{p_n(\mathbf{x}|\mathbf{y}) + \tilde{p}_\theta(\mathbf{x}|\mathbf{y})} \right] + \mathbb{E}_{p_{\text{data}}(\mathbf{x}|\mathbf{y})} \left[\log \frac{\tilde{p}_\theta(\mathbf{x}|\mathbf{y})}{p_n(\mathbf{x}|\mathbf{y}) + \tilde{p}_\theta(\mathbf{x}|\mathbf{y})} \right] \\
&= \mathbb{E}_{p_n(\mathbf{x}|\mathbf{y})} \left[\log \frac{1}{1 + p_\theta(\mathbf{x}|\mathbf{y})} \right] + \mathbb{E}_{p_{\text{data}}(\mathbf{x}|\mathbf{y})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{y})}{1 + p_\theta(\mathbf{x}|\mathbf{y})} \right] \\
&= \mathbb{E}_{p_n(\mathbf{x}|\mathbf{y})} \left[\log \frac{\exp(-f_{\mathbf{x}}(\mathbf{x}, \mathbf{y}))}{\exp(-f_{\mathbf{x}}(\mathbf{x}, \mathbf{y})) + 1} \right] + \mathbb{E}_{p_{\text{data}}(\mathbf{x}|\mathbf{y})} \left[\log \frac{1}{\exp(-f_{\mathbf{x}}(\mathbf{x}, \mathbf{y})) + 1} \right] \\
&= \mathbb{E}_{p_n(\mathbf{x}|\mathbf{y})} \left[\log (1 - \sigma(f_{\mathbf{x}}(\mathbf{x}, \mathbf{y}))) \right] + \mathbb{E}_{p_{\text{data}}(\mathbf{x}|\mathbf{y})} \left[\log \sigma(f_{\mathbf{x}}(\mathbf{x}, \mathbf{y})) \right].
\end{aligned} \tag{22}$$

Thus, the final EBM-NCE contrastive SSL objective is

$$\begin{aligned}
\mathcal{L}_{\text{EBM-NCE}} = & -\frac{1}{2} \mathbb{E}_{p_{\text{data}}(\mathbf{y})} \left[\mathbb{E}_{p_n(\mathbf{x}|\mathbf{y})} \log (1 - \sigma(f_{\mathbf{x}}(\mathbf{x}, \mathbf{y}))) + \mathbb{E}_{p_{\text{data}}(\mathbf{x}|\mathbf{y})} \log \sigma(f_{\mathbf{x}}(\mathbf{x}, \mathbf{y})) \right] \\
& -\frac{1}{2} \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \left[\mathbb{E}_{p_n(\mathbf{y}|\mathbf{x})} \log (1 - \sigma(f_{\mathbf{y}}(\mathbf{y}, \mathbf{x}))) + \mathbb{E}_{p_{\text{data}}(\mathbf{y}|\mathbf{x})} \log \sigma(f_{\mathbf{y}}(\mathbf{y}, \mathbf{x})) \right].
\end{aligned} \tag{23}$$

E IMPLEMENTATION AND EXPERIMENT DETAILS

Editing function. For the editing function, we consider both linear (Eqs. (24) and (25)) and non-linear (Eq. (26)) cases as below, *i.e.*, for $\tilde{\mathbf{z}}_{i,\alpha} = h(\mathbf{z}, \mathbf{d}_i, \alpha)$:

$$\tilde{\mathbf{z}}_{i,\alpha} = \mathbf{z} + \alpha \cdot \mathbf{d}_i, \quad \mathbf{d}_i = \text{norm} \circ \text{Linear}(\mathbf{e}_i), \quad (24)$$

$$\tilde{\mathbf{z}}_{i,\alpha} = \mathbf{z} + \alpha \cdot \mathbf{d}_i, \quad \mathbf{d}_i = \text{sqrt} \circ \text{norm} \circ \text{ReLU} \circ \text{Linear}(\mathbf{e}_i), \quad (25)$$

$$\tilde{\mathbf{z}}_{i,\alpha} = \mathbf{z} + \alpha \cdot \mathbf{d}_i + \text{norm} \circ \text{Linear} \circ \text{ReLU} \circ \text{Linear}(\mathbf{z} \oplus \mathbf{d}_i \oplus [\alpha]), \quad \mathbf{d}_i = \text{norm} \circ \text{Linear} \circ \text{ReLU} \circ \text{Linear}(\mathbf{e}_i), \quad (26)$$

where \circ means the composition of two functions.

Objective function. The objective function is given as:

$$\mathcal{L} = c_1 \cdot \mathbb{E}_{u,v}[\mathcal{L}_{\text{GraphCG-NCE}}] + c_2 \cdot \mathcal{L}_{\text{sim}} + c_3 \cdot \mathcal{L}_{\text{sparsity}}, \quad (27)$$

where $\mathcal{L}_{\text{GraphCG-NCE}}$ is the MI estimation defined in Eq. (9), \mathcal{L}_{sim} is the direction similarity defined in Eq. (10), and $\mathcal{L}_{\text{sparsity}}$. c_1 , c_2 and c_3 are three coefficients accordingly.

Hyperparameters. We list the key hyperparameters in Table 4, and all the results are evaluated on 100 sampled sequences. We also want to highlight that DisCo is an unstable baseline, in the sense that once we add more training data (*e.g.*, from 100 to 500) or more training epochs (*e.g.*, from 1 epoch to 5 epochs), the model will collapse, with the nan loss. Thus, here we are reporting the most reasonable results for DisCo, *i.e.*, 100 training data with 1 epoch.

Table 4: Hyperparameter specifications.

	Hyperparameter	Value
Random	D	{10}
	α	{-3, -2.7, -2.4, ..., 2.7, 3}
Variance	D	{10}
	α	{-3, -2.7, -2.4, ..., 2.7, 3}
	# training data	{100, 500}
SeFa	D	{10}
	α	{-3, -2.7, -2.4, ..., 2.7, 3}
	# training data	{100, 500}
DisCo	D	{10}
	α	{-3, -2.7, -2.4, ..., 2.7, 3}
	# training data	{100}
	epochs	{1}
GraphCG	D	{10}
	α	{-3, -2.7, -2.4, ..., 2.7, 3}
	# training data	{100, 500}
	epochs	{20, 100}
	coefficient c_1	{0, 1}
	coefficient c_2	{0, 1}
	coefficient c_3	{0, 1}

Hardware. We use V100 GPU cards, and each job (w.r.t. different hyperparameters) for GraphCG can be finished within 3 hours on a single GPU card.

Time complexity. The time complexity of GraphCG is $O(B \times D^2)$ for GraphCG-P and $O(B^2 \times D^2)$ for GraphCG-R, where B is the number of data points for each batch. Here we omit the constants for step-sizes.

F RESULTS: MOLECULAR GRAPH

F.1 EVALUATION METRICS

Change of Structure Factors and Calibrated Tanimoto Similarity (CTS). We are interested in the change of the graph structure (the steerable factors) along the output sequence edited with the i -th direction. To evaluate the structure change, we apply the Tanimoto similarity between each output molecule and the anchor molecule, as shown in Fig. 5(a). Besides, for the ease of evaluating the monotonicity, we utilize a transformation (on the Tanimoto similarity) of output molecules with positive step size by taking the deduction from 2. We call this calibrated Tanimoto similarity sequence (CTS), *i.e.*, $\{s(\bar{x}')\}_i$, as shown in Fig. 5(b).

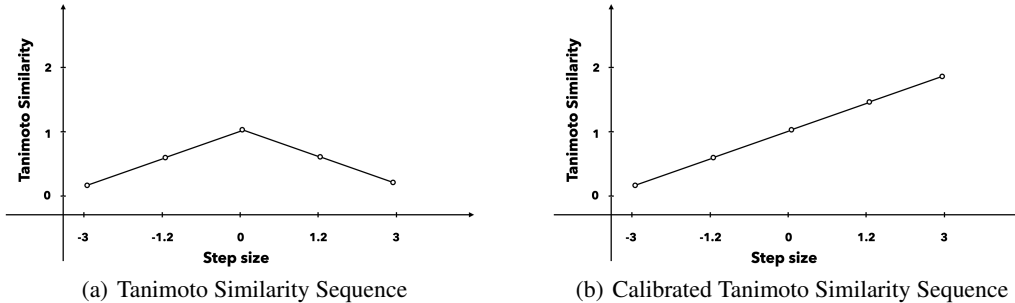


Figure 5: Fig. 5(a) is the original Tanimoto similarity sequence w.r.t. the anchor molecule, *i.e.*, step size with 0 in the figure. Yet, this is not easy to compute the monotonicity. We thus propose the calibrated Tanimoto similarity sequence, by taking the deduction from 2 for output molecules with positive step size, as shown in Fig. 5(b).

Sequence Monotonic Ratio (SMR). For evaluation, we propose a metric called Sequence Monotonic Ratio (SMR), $\phi_{\text{SMR}}(\gamma, \tau)_i$. It measures the monotonic ratio of M generated sequences edited with the i -th direction. It has two arguments: the diversity threshold γ constrains the minimum number of distinct molecules, and the tolerance threshold τ controls the non-monotonic tolerance ratio along each sequence.

In specific, for each learned semantic direction i , we will generate M sequences of edited molecules, and the calibrated Tanimoto similarity for each sequence is marked as $\{s(\bar{x}')\}_i^m$. Then we can define the SMR on each direction as:

$$\begin{aligned} \phi_{\text{SMR}}(\gamma, \tau)_i &= \frac{1}{M} \sum_{m=1}^M \phi_{\text{SMR}}(\{s(\bar{x}')\}_i^m, \gamma, \tau), \\ \phi_{\text{SMR}}(\{s(\bar{x}')\}_i^m, \gamma, \tau) &= \begin{cases} 1, & \text{len}(\text{set}\{s(\bar{x}')\}_i^m) \geq \gamma \wedge \text{monotonic}_{\tau}(\{s(\bar{x}')\}_i^m) \\ 0, & \text{otherwise} \end{cases} \end{aligned} \quad (28)$$

Evaluating the Diversity of Semantic Directions. SMR can evaluate all the output sequences generated by one direction. To better illustrate that GraphCG is able to learn multiple directions with various semantic information, we also consider taking the average of top- K SMR w.r.t. directions to reveal that all the best K directions are semantically meaningful, as in Eq. (29):

$$\text{top-K}(\gamma, \tau) = \frac{1}{K} \sum_{i \in \text{top-K directions}} (\phi_{\text{SMR}}(\gamma, \tau)_i). \quad (29)$$

F.2 RESULTS ON MOLECULAR STRUCTURES

Next we would like to show the comprehensive SMR on the CTS results with respect to different backbone models, as in Tables 5 and 6.

Table 5: This table lists the sequence monotonic ratio (SMR, %) on calibrated Tanimoto similarity (CST) w.r.t. the top-1, top-2, and top-3 directions. The backbone model is the pretrained MoFlow on ZINC250k.

Edit Method	top-K γ τ	Tanimoto top-1						Tanimoto top-2						Tanimoto top-3					
		2		3		4		2		3		4		2		3		4	
		0	0.2	0	0.2	0	0.2	0	0.2	0	0.2	0	0.2	0	0.2	0	0.2	0	0.2
Random		35.0	36.0	23.0	25.0	12.0	15.0	34.5	36.0	22.5	25.0	11.5	14.0	34.0	36.0	22.0	24.0	11.0	13.7
Variance		32.0	36.0	24.0	28.0	12.0	16.0	31.5	35.5	21.0	26.5	10.5	16.0	30.3	35.3	20.0	25.0	10.0	15.0
SeFa		23.0	23.0	4.0	4.0	0.0	0.0	19.0	19.0	4.0	4.0	0.0	0.0	17.3	17.3	3.3	3.3	0.0	0.0
DisCo		8.0	15.0	7.0	14.0	2.0	8.0	7.5	13.5	6.0	12.5	2.0	8.0	7.0	13.0	5.3	11.7	2.0	7.7
GraphCG-P	Eq. (24)	39.0	40.0	27.0	28.0	15.0	18.0	38.5	40.0	26.0	28.0	15.0	18.0	37.0	39.0	24.7	27.3	14.3	17.3
	Eq. (25)	35.0	37.0	19.0	22.0	8.0	11.0	33.5	36.5	18.5	21.5	7.0	10.5	31.7	34.7	17.7	20.7	6.3	9.3
	Eq. (26)	44.0	46.0	32.0	34.0	16.0	18.0	42.5	44.5	30.0	32.0	15.0	17.5	41.7	44.0	29.0	31.0	13.7	16.3
GraphCG-R	Eq. (24)	37.0	42.0	25.0	26.0	11.0	14.0	37.0	40.0	23.0	25.5	11.0	13.5	36.3	39.0	22.0	24.3	10.3	13.3
	Eq. (25)	35.0	37.0	19.0	22.0	8.0	11.0	33.5	36.5	18.5	21.5	7.0	10.5	31.7	34.7	17.7	20.7	6.3	9.3
	Eq. (26)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Table 6: This table lists the sequence monotonic ratio (SMR, %) on calibrated Tanimoto similarity (CST) w.r.t. the top-1, top-2, and top-3 directions. The backbone model is the pretrained HierVAE on ChEMBL.

Edit Method	top-K γ τ	Tanimoto top-1						Tanimoto top-2						Tanimoto top-3					
		2		3		4		2		3		4		2		3		4	
		0	0.2	0	0.2	0	0.2	0	0.2	0	0.2	0	0.2	0	0.2	0	0.2	0	0.2
Random		14.0	45.0	14.0	45.0	14.0	43.0	11.0	43.5	11.0	43.5	11.0	42.5	10.0	42.3	10.0	42.3	9.3	41.7
Variance		28.0	64.0	23.0	59.0	19.0	55.0	22.5	59.5	19.5	57.0	17.5	55.0	20.3	54.3	18.3	52.7	15.7	50.3
SeFa		4.0	41.0	4.0	41.0	4.0	41.0	3.0	41.0	3.0	41.0	3.0	41.0	2.3	36.0	2.3	36.0	2.3	36.0
GraphCG-P	Eq. (24)	23.0	61.0	19.0	57.0	15.0	53.0	19.0	59.0	16.5	56.5	13.5	53.0	17.0	55.3	15.3	53.7	12.7	50.7
	Eq. (25)	62.0	77.0	40.0	73.0	32.0	65.0	60.5	74.0	38.5	67.5	29.0	61.0	59.3	70.3	36.0	64.3	26.3	57.7
	Eq. (26)	29.0	71.0	28.0	70.0	27.0	69.0	22.0	62.0	21.5	61.5	20.5	61.0	18.7	57.7	18.3	57.3	17.7	56.7
GraphCG-R	Eq. (24)	16.0	56.0	16.0	56.0	15.0	55.0	13.5	48.0	13.5	48.0	12.0	47.0	11.7	44.7	11.7	44.7	10.7	43.3
	Eq. (25)	61.0	74.0	42.0	67.0	30.0	55.0	59.0	69.5	40.0	64.0	29.5	54.5	57.7	67.7	38.0	62.3	28.7	53.7
	Eq. (26)	25.0	57.0	24.0	57.0	21.0	55.0	20.0	55.0	19.5	54.5	17.0	52.0	17.3	52.7	17.0	52.3	15.0	50.0

F.3 VISUALIZATION

For a more comprehensive visualization of the steerable factors in molecular graphs, we demonstrate 16 molecular graph paths along the 4 selected directions in Fig. 6, and the backbone DGM is HierVAE pretrained on ChEMBL. The CTS holds good monotonic trend in all these sequences. Each direction shows certain unique changes in the molecular structures, *i.e.*, the steerable factors in molecules. Some structural changes are reflected in molecular properties. We expand all the details below. In Fig. 6(a) and Fig. 6(b), the number of halogen atoms and hydroxyl groups (in alcohols and phenols) in the molecules decrease from left to right, respectively. In Fig. 6(c), the number of amides in the molecules increases along the path. As a result, the topological polar surface area (tPSA) of the molecules increase accordingly, which is a key molecular property for the prediction of drug transport properties, *e.g.*, permeability [12]. In Fig. 3(d), the flexible chain length, marked by the number of ethylene (CH_2CH_2) units, increases from left to right. Since the number of rotatable bonds (NRB) measures the molecular flexibility, it also increases accordingly [57].

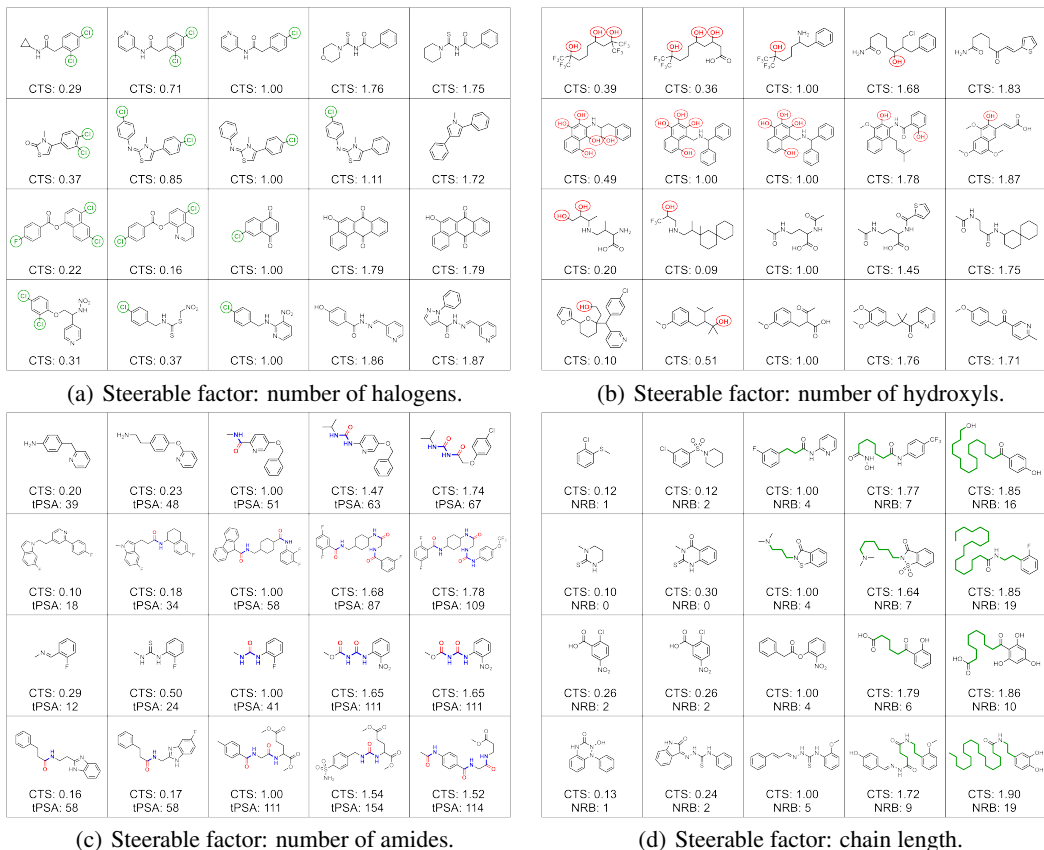


Figure 6: GraphCG for molecular graph editing. We visualize the output molecules and CTS on four directions with four sequences each, where each sequence consists of five steps. The center point is the anchor molecule, and the other four points correspond to step size with -3, -1.8, 1.8, and 3 respectively. Fig. 6(a) to Fig. 6(c) show how functional groups in the molecules can be viewed as the steerable factors as they change along the path, such as halogen atoms, hydroxyl groups and amides. Fig. 6(d) illustrates the effect on the steerable factor on the length of flexible chains in the molecules. Notably, certain properties change together with molecular structures, like topological polar surface area (tPSA) and number of rotatable bonds (NRB).

Since we have determined four directions with semantic information matching with the domain knowledge (fragments), then we can check if the disentanglement measure changes before and after editing. Notice that comparing to Table 1, three measures are infeasible (*i.e.*, with value nan), and we show three feasible ones in Table 7.

Table 7: Disentanglement measure before and after editing. The corresponding model is the GraphCG-P with Eq. (25), and the backbone generative model is HierVAE pretrained on ChEMBL. The better performance is marked in **bold**.

Fragment		BetaVAE \uparrow	MIG \uparrow	SAP \uparrow
Halogen	after editing	0.617	0.010	0.017
	before editing	0.950	0.062	0.017
Hydroxyls	after editing	0.833	0.031	0.017
	before editing	0.933	0.113	0.067
Amide	after editing	0.400	0.041	0.017
	before editing	0.933	0.136	0.017
Chain length	after editing	0.400	0.051	0.000
	before editing	0.700	0.020	0.017

G RESULTS: POINT CLOUDS

Here we compare two editing functions, *i.e.*, the key function design in GraphCG. We provide the visualizations for the linear editing function in Eq. (25), and non-linear editing function in Eq. (26).

First we want to highlight that all the samples are generated randomly. In the linear case in Appendix G.1, we can observe that the shape of the airplanes, cars, and chairs, are steerable using GraphCG. We also find it interesting that GraphCG can steer more finger-trained factors, like modifying the airplane engines. However, in the non-linear case, the diversity of the edited data is smaller. This can be observed from the middle columns in Appendix G.2. We will leave this for future exploration.

G.1 LINEAR EDITING FUNCTION



Figure 7: GraphCG for point clouds (Airplane) editing. It can successfully reflect these steerable factors: engine, fuselage length, wing size, wing shape, and wing thickness.

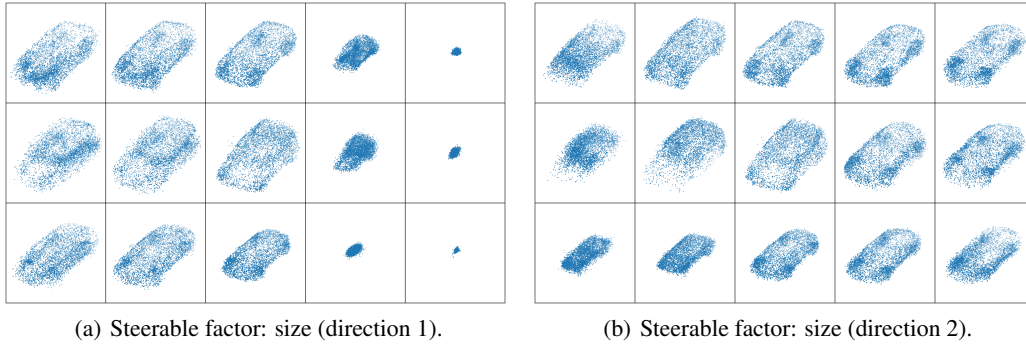


Figure 8: GraphCG for point clouds (Car) editing. The steerable factors on this dataset are not obvious, and here we only plot the car size editable with two directions.

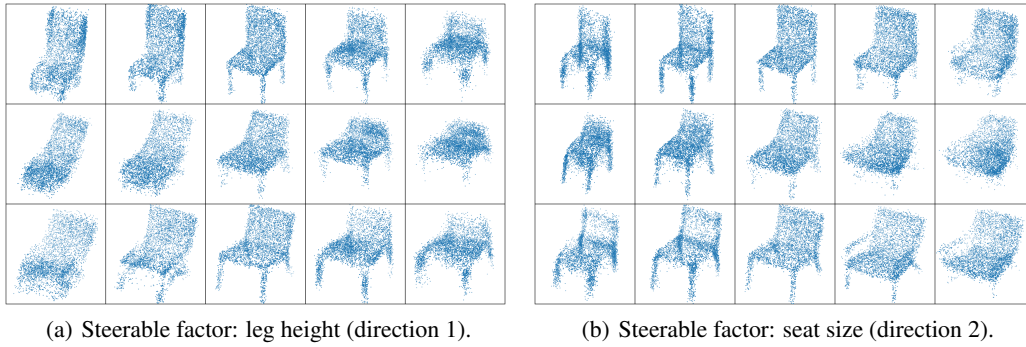


Figure 9: GraphCG for point clouds (Chair) editing. It can successfully reflect these steerable factors: leg height, and seat size.

G.2 NON-LINEAR EDITING FUNCTION

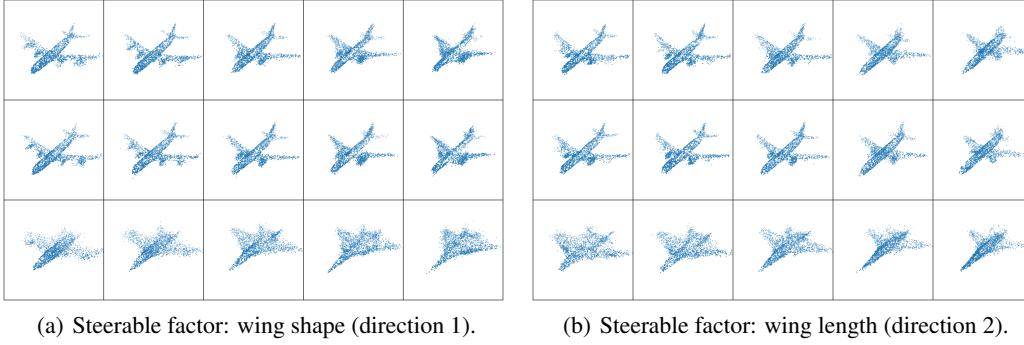


Figure 10: GraphCG for point clouds (Airplane) editing. It can successfully reflect these steerable factors: wing shape and wing length.

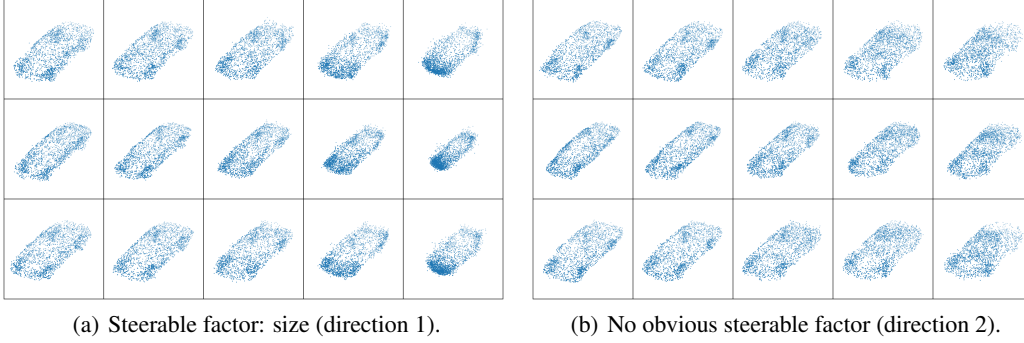


Figure 11: GraphCG for point clouds (Car) editing. The steerable factors on this dataset are not obvious, and here we only plot the car size editable with one directions.

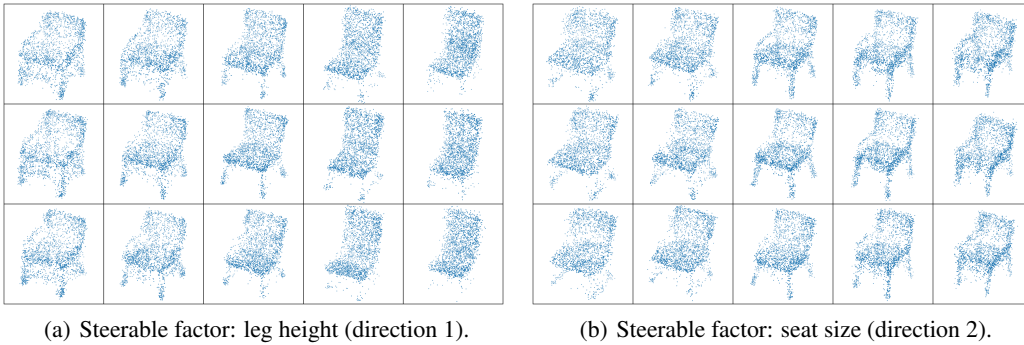


Figure 12: GraphCG for point clouds (Chair) editing. It can successfully reflect these steerable factors: leg height, and seat size.