

---

# Order Agnostic Autoregressive Graph Generation

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Graph generation is a fundamental problem in various domains, including chemistry  
2 and social networks. Recent work has shown that molecular graph generation  
3 using recurrent neural networks (RNNs) is advantageous compared to traditional  
4 generative approaches which require converting continuous latent representations  
5 into graphs. One issue which arises when treating graph generation as sequential  
6 generation is the arbitrary order of the sequence which results from a particular  
7 choice of graph flattening method. In this work we propose using RNNs, taking into  
8 account the non-sequential nature of graphs by adding an Orderless Regularization  
9 (OLR) term that encourages the hidden state of the recurrent model to be invariant  
10 to different valid orderings present under the training distribution. We demonstrate  
11 that sequential graph generation models benefit from our proposed regularization  
12 scheme, especially when data is scarce. Our findings contribute to the growing body  
13 of research on graph generation and provide a valuable tool for various applications  
14 requiring the synthesis of realistic and diverse graph structures.

## 15 1 Introduction

16 Graphs are powerful representations of complex relationships and structures found in a wide range  
17 of domains, including social networks, molecular chemistry, transportation networks, distributed  
18 algorithms and many more. A dedicated class of architectures, Graph Neural Networks (GNNs), has  
19 been developed to handle the specific properties of graphs. Graphs are naturally versatile objects,  
20 but such versatility comes at the cost of lack of structure and no naturally induced order. Most GNN  
21 architectures therefore operate by applying a neural architecture at the node level followed by an  
22 aggregation step which takes into account the local neighborhood structure of the graph. By stacking  
23 multiple such layers, a GNN is able to perform node-level or graph-level tasks that take into account  
24 the entire structure of the graph.

25 The ability to generate realistic and structured graphs is essential for various applications ranging from  
26 drug design [17, 20, 34, 36, 44, 55, 66, 69] to program synthesis [9, 32, 11, 1, 30, 65, 7, 15]. In recent  
27 years a wide variety of generative models have been developed, including generative adversarial  
28 networks (GANs), variational autoencoders (VAEs), normalizing flows, and diffusion models. These  
29 algorithms devise different strategies to learn continuous mappings from a latent distribution to a space  
30 of realistic examples. Unfortunately, graphs do not admit a natural representation in a continuous  
31 space; consequently, the discrete and unordered nature of graphs make them less amenable to the  
32 methods mentioned above for the task of graph generation. A different type of generative model  
33 relies on autoregressive architectures which enable processing a sequence and generating the next  
34 element; for example, these architectures are commonly used for large language models. Generally,  
35 autoregressive models are applicable when the generated objects admit a sequential order.

36 In this work we focus on sequential generation of graphs using autoregressive neural architectures. A  
37 strong motivating factor for choosing autoregressive architectures is that we are particularly interested  
38 in molecular graph generation; and in this context Flam *et al.* [23] have shown that sequential

39 generation is favorable compared to other approaches. The specific representation we consider is  
40 depth-first search (DFS) trajectories of graphs. The reasons for this choice of representation is  
41 twofold: (a) DFS is a natural way of flattening graphs into sequences; (b) in the chemistry community  
42 DFS is used to convert molecules into strings. However, an issue arises when converting graphs into  
43 sequences: there are many DFS trajectories for a given graph. Indeed, for many graph flattening  
44 methods, there is an arbitrariness in the order of the sequence which results [63, 12].

45 In order to alleviate the dependency on a specific trajectory, we add a regularization term dubbed  
46 Orderless Regularization (OLR) which ensures the learnt model is invariant to different DFS orderings  
47 of the same graph. For the sake of training with OLR, one needs to generate different DFS trajectories  
48 with a common end-vertex which is known to be hard [5]. We formalize the notion of graph-level  
49 invariance and devise an efficient algorithm to generate such trajectories under certain constraints.  
50 Finally, we demonstrate empirically that our regularization term is beneficial when the amount of  
51 training data is limited by considering the use case of small molecule generation.

52 The remainder of the paper is structured as follows: in Section 2 we provide background and introduce  
53 the concepts, definitions, and notations used throughout the paper. Section 3 goes into the details of  
54 OLR over DFS trajectories. Section 4 is devoted to related work. In Section 5 we provide empirical  
55 evidence for the effectiveness of OLR. Section 6 provides concluding remarks.

## 56 2 Background

57 In this section we formally define the problem of graph generation, the notations and definitions  
58 necessary to present our proposed method. We denote matrices by bold uppercase letters,  $\mathbf{M} \in \mathbb{R}^{n \times m}$ ,  
59 vectors by bold lowercase letters,  $\mathbf{v}$ , and the  $i^{\text{th}}$  entry of  $\mathbf{v}$  by  $v_i$ . We proceed with a general  
60 formulation of recurrent models.

### 61 2.1 Recurrent Models

62 Let  $\mathcal{X}$ ,  $\mathcal{H}$ , and  $\mathcal{Y}$  be the spaces of inputs, hidden states, and outputs, respectively. Given an input  
63 sequence  $\mathbf{x} \equiv (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathcal{X}^n$ , a recurrent model consists of two functions, the state update  
64 function  $f_u : \mathcal{H} \times \mathcal{X} \rightarrow \mathcal{H}$ ,

$$\mathbf{h}_{t+1} = f_u(\mathbf{h}_t, \mathbf{x}_t), \quad (2.1)$$

65 and the output function  $f_o : \mathcal{H} \times \mathcal{X} \rightarrow \mathcal{Y}$ ,

$$\mathbf{y}_t = f_o(\mathbf{h}_t, \mathbf{x}_t). \quad (2.2)$$

66 where  $\mathbf{h}_0 \in \mathcal{H}$ . We overload the notation and denote the hidden state and output of a recurrent model  
67 over a sequence as  $f_u(\mathbf{x})$  and  $f_o(\mathbf{x})$  respectively.

68 The formulation presented of recurrent models is broad and able to capture RNNs as well as more  
69 complex architectures such as Gated Recurrent Units (GRUs, [13]) and Long-Short Term Memory  
70 networks (LSTMs, [33]). For example, in the case of a very simple, vanilla RNN,<sup>1</sup>

$$\mathbf{h}_{t+1} = \sigma_h(\mathbf{A}\mathbf{h}_t + \mathbf{B}\mathbf{x}_t) \quad (2.3)$$

71 and

$$\mathbf{y}_t = \sigma_y(\mathbf{C}\mathbf{h}_t + \mathbf{D}\mathbf{x}_t) \quad (2.4)$$

72 where  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$  are matrices with the appropriate dimensions and  $\sigma_h, \sigma_y$  are standard non-  
73 linearities such as *sigmoid* or *tanh*.

### 74 2.2 Graph Generation

75 A graph is given by  $G = (\mathcal{V}, \mathcal{E})$  where  $\mathcal{V}$  is a set of nodes (or vertices) and  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$  is a set of  
76 tuples denoting the nodes connected by an edge in the graph. Additionally, for each  $v \in \mathcal{V}$ , denote  
77 by  $x_v \in \mathbb{R}^m$  the features of node  $v$ . Similarly,  $e_{uv} \in \mathbb{R}^k$  denotes the features of the edge  $(u, v) \in \mathcal{E}$ .  
78 For example, in a molecular graph, nodes are atoms, and their features will contain the element; and  
79 edges correspond to bonds, and their features will contain the bond types (single, double, etc.).<sup>2</sup>

<sup>1</sup>Note that the bias term may be encapsulated into the input processing matrices by expanding the input with an additional dimension and assigning a fixed value of 1 on that coordinate.

<sup>2</sup>Molecular node and edge features may contain other properties as well.

80 Another example is of social networks, where nodes corresponds to users and their features to user  
81 profiles; and edges correspond to connections between users and their features contain metadata on  
82 this connection.

83 The topic of designing neural networks to operate specifically on graphs is dominated by Graph  
84 Neural Networks (GNNs) which mostly rely on a message-passing scheme to propagate information  
85 between nodes. While these architectures are extremely successful in node level and graph level  
86 prediction, they are not as prevalent in the context of graph generation, and many such approaches  
87 are restricted to small graphs (though [16] is a notable exception).

88 Formally, the task of graph generation is usually concerned with learning to model distributions:  
89 concretely, given a set of  $N$  graphs  $\{G_i\}_{i=1}^N$  originating from an underlying distribution  $p$ , the  
90 goal of graph generation is to devise an algorithm that generates new graphs from the underlying  
91 distribution  $p$ . Prior work has mostly adapted successful generative methods over a continuous space  
92 to the domain of graphs [25, 8, 37, 51]. In this work we focus on using recurrent models which  
93 can be employed naturally to generate discrete objects. Crucially, Flam *et al.* [23] have shown that  
94 sequential generation is favorable compared to competing approaches in the context of molecular  
95 graph generation.

## 96 2.3 Sequential Graph Generation

97 When applying recurrent models for graph generation, the graph first needs to be “flattened” into  
98 a sequence. As there is no natural order for a graph, one must artificially induce such an order; for  
99 example, the approach taken by [67] considers generation of breadth-first search (BFS) trajectories.  
100 While there are many ways to convert a graph into a sequence, in this work we focus on depth-first  
101 search (DFS); a strong motivation for this choice is that this is the method used to convert graph  
102 molecules into a linear representation called SMILES strings [64]. By convention, the output of the  
103 DFS algorithm is a spanning tree and we consider the induced order of the graph as the order in  
104 which the vertices were visited during the DFS run (also known as pre-order traversal).

105 In what follows we formally define the concepts discussed.

106 **Definition 2.1.** Given a connected graph  $G = (\mathcal{V}, \mathcal{E})$  with  $|\mathcal{V}| = n$ , we say the permutation  $\pi \in \mathbb{S}_n$   
107 is a valid ordering of  $G$  if it is possible to run DFS over  $G$  and visit the vertices in the order induced  
108 by  $\pi$ . Denote the sequence corresponding to a valid ordering  $\pi$  of  $G$  by

$$\mathbf{s}(G, \pi) = (v_{\pi(1)}, \dots, v_{\pi(n)}). \quad (2.5)$$

109 Denote the set of all such sequences for a given graph  $G$  by

$$\mathcal{S}(G) = \{\mathbf{s}(G, \pi) : \pi \text{ is a valid ordering of } G\}. \quad (2.6)$$

110 Clearly, for a non-trivial graph  $\mathcal{S}(G)$  will contain many sequences. In this work we have a special  
111 interest in sequences that share the same end vertex.

112 **Definition 2.2.** Let  $\mathcal{S}(G, v)$  denote all sequences terminating at node  $v \in \mathcal{V}$ , formally,

$$\mathcal{S}(G, v) = \{\mathbf{s} \in \mathcal{S}(G) : s_n = v\} \quad (2.7)$$

113 In the following section we discuss the desired properties for recurrent models when used for graph  
114 generation.

## 115 3 Structure Agnostic Recurrent Models

116 Recurrent models are a natural choice when generating discrete objects such as text. On the other  
117 hand, graphs are discrete objects with no naturally induced order. In Section 2 we described a  
118 mapping between graphs and sequences; and in particular, the fact that many different sequences  
119 correspond to the same graph. In this section we present our method that overcomes the issues  
120 described.

### 121 3.1 Generating Depth-First Search Traversals

122 In this work we use recurrent models to generate DFS traversals of graphs. Clearly when generating a  
123 DFS traversal, the next node to be generated depends on the nodes generated thus far and in particular

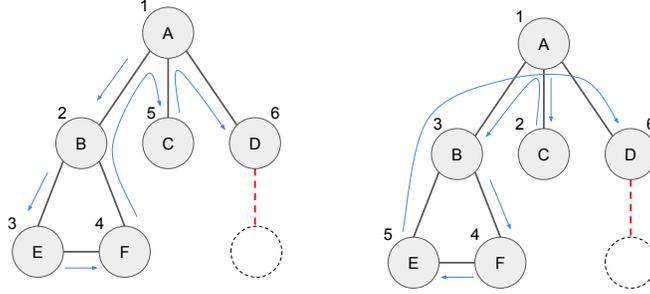


Figure 1: Illustration of two DFS traversals of the same graph starting from node  $A$  and terminating at node  $D$ , blue lines denote traversal order. (Left) traversal resulting in the sequence  $A(BEF)(C)D$ . (Right) traversal resulting in the sequence  $A(C)(BFE)D$ . The parentheses denote the opening and closing of branches when traversing the tree; with this syntax it is possible to reconstruct the tree from such sequences. Note that multiple sequences correspond to the same tree, a fact that lies at the heart of this work.

124 the last generated node. An important observation is that the output of the recurrent model should be  
 125 *invariant to different valid orderings corresponding to the same subgraph* as long as they lead to the  
 126 same node. The following definition formalizes this notion,

127 **Definition 3.1.** We say a recurrent model is **structure invariant** with respect to a connected graph  $G$   
 128 if

$$\forall v \in \mathcal{V}, \forall \mathbf{s}_1, \mathbf{s}_2 \in \mathcal{S}(G, v) \text{ it is the case that } f_o(\mathbf{s}_1) = f_o(\mathbf{s}_2). \quad (3.1)$$

129 If the above condition is satisfied for all  $G \sim \mathcal{D}$ , we say that the recurrent model is **structure**  
 130 **invariant with respect to a distribution**  $\mathcal{D}$ .

131 Figure 1 depicts a graph and two different DFS traversals sharing the same root and terminal node. A  
 132 recurrent model processing the two DFS traversals will ideally generate the same node that will be  
 133 attached to node  $D$ .

134 Definition 3.1 describes the structure invariance property with respect to a graph. Since recurrent  
 135 models generate the traversal sequentially, we would like this property to hold at any moment during  
 136 generation, i.e., we want to modify Definition 3.1 to take into account partial DFS traversals.

137 **Definition 3.2.** For a connected graph  $G$ , we say a connected subgraph  $\tilde{G} \subseteq G$  is **induced by DFS**  
 138 **over**  $G$  if there exists a valid ordering  $\pi \in S_n$  of  $G$ , and  $k \leq n$  such that  $(v_{\pi(1)}, \dots, v_{\pi(k)})$  is a valid  
 139 ordering of  $\tilde{G}$ . Denote the set of all DFS induced subgraphs over  $G$  by  $\mathcal{G}_{DFS}(G)$ .

140 At this stage, a reader might question the necessity of Definition 3.2 and why  $\mathcal{G}_{DFS}(G)$  differs from  
 141 the set of all connected subgraphs of  $G$ . We note that for a general connected graph  $\mathcal{G}_{DFS}(G)$  does  
 142 **not** correspond to the set of all connected subgraphs.

143 **Proposition 3.3.** For a connected graph  $G$ ,

$$\mathcal{G}_{DFS}(G) \neq \left\{ \tilde{G} \mid \tilde{G} \subseteq G \text{ and } \tilde{G} \text{ is connected} \right\} \quad (3.2)$$

144 Figure 2 depicts a graph and two connected subgraphs, one which is induced by DFS and the other  
 145 that cannot be obtained by a DFS traversal.

146 With the notion of DFS induced subgraphs in hand, we now present the following definition of *total*  
 147 *structure invariance*:

148 **Definition 3.4.** We say a recurrent model is **totally structure invariant** with respect to a connected  
 149 graph  $G$ , if

$$\forall \tilde{G} \in \mathcal{G}_{DFS}(G), \forall v \in \mathcal{V}(\tilde{G}), \forall \mathbf{s}_1, \mathbf{s}_2 \in \mathcal{S}(\tilde{G}, v) \text{ it is the case that } f_o(\mathbf{s}_1) = f_o(\mathbf{s}_2). \quad (3.3)$$

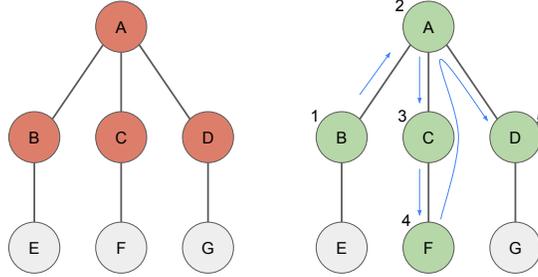


Figure 2: Illustration of the same graph with two connected subgraphs: (Left) subgraph which is not induced by DFS. (Right) subgraph induced by DFS, arrows depict a traversal resulting in the sequence  $BA(CF)D$ .

150 *If the above condition is satisfied for all  $G \sim \mathcal{D}$ , we say that the recurrent model is **totally structure***  
 151 ***invariant with respect to a distribution  $\mathcal{D}$ .***

152 In the next section we discuss how to train recurrent models which are totally structure invariant with  
 153 respect to a given training distribution over graphs.

### 154 3.2 Regularizing Towards Total Structure Invariance

155 Motivated by the observation discussed in Section 3.1, we propose training recurrent models that  
 156 are totally structure invariant with respect to the underlying distribution over graphs. It would be  
 157 appealing to characterize the class of all totally structure invariant functions and optimize over those.  
 158 Unfortunately, it is difficult to attain a crisp characterization of structure invariance as this property  
 159 depends on the training distribution.

160 Instead, we propose encouraging total structure invariance via regularization. Specifically, we would  
 161 like to minimize the following auxiliary loss,

$$\mathbb{E}_{G \sim \mathcal{D}} \mathbb{E}_{\tilde{G} \sim \mathcal{G}_{DFS}(G)} \mathbb{E}_{v \in \mathcal{V}(\tilde{G})} \mathbb{E}_{\mathbf{s}_1, \mathbf{s}_2 \in \mathcal{S}(\tilde{G}, v)} \left[ (f_o(\mathbf{s}_1) - f_o(\mathbf{s}_2))^2 \right] \quad (3.4)$$

162 which we refer to as *Orderless Regularization (OLR)*. Examining Equation 3.4, we note that sampling  
 163 from  $\mathcal{G}_{DFS}(G)$  is easily done by randomly selecting a root node and running DFS with stochastic  
 164 decision making. On the other hand, given  $\tilde{G}$  and  $v$ , sampling from  $\mathcal{S}(\tilde{G}, v)$  is hard and has been  
 165 shown to be NP-complete [5].

### 166 3.3 Sampling Trajectories with Common End Vertex

167 The problem of generating all DFS trajectories that terminate at the same vertex is hard and there are  
 168 no known efficient algorithms for this task. In order to overcome this obstacle we apply a heuristic  
 169 for computing such trajectories. We highlight that our proposed scheme is not equivalent to a uniform  
 170 sampling over all possible trajectories; however, in Section 5 we show that the resulting regularization  
 171 scheme is effective empirically.

172 Next, we formally show that for practical graphs there exists efficient algorithms to generate such  
 173 trajectories.

174 **Definition 3.5.** *Let  $G = (\mathcal{V}, \mathcal{E})$  be an arbitrary graph.  $G$  is said to be  **$k$ -edge-connected** if the*  
 175 *subgraph  $G' = (\mathcal{V}, \mathcal{E} \setminus \tilde{\mathcal{E}})$  is connected for all  $\tilde{\mathcal{E}} \subseteq \mathcal{E}$  such that  $|\tilde{\mathcal{E}}| < k$ .*

176 **Proposition 3.6.** *There is an efficient algorithm to find distinct DFS trajectories with common end*  
 177 *vertex for any  $k$ -edge connected graph for  $k \leq 2$ .*

178 We note that in many real world tasks, graph are rarely  $k$ -edge connected for  $k > 2$ . For example, in  
179 the ZINC molecular dataset, more than 99.5% of molecular graphs are 1-edge connected.

180 *Proof Sketch.* Find a min-cut: by the definition of 1-edge-connected graphs the min-cut includes  
181 a single crossing edge  $(u, v)$ . By removing  $(u, v)$  the graph is partitioned into two connected  
182 components,  $G_1$  and  $G_2$  containing  $u$  and  $v$  respectively. Run a DFS on  $G_1$  with  $u$  as the root vertex  
183 to result in  $(u_1, \dots, u_k)$ , and similarly for  $G_2$  to result in  $(v_1, \dots, v_m)$  (where  $v_1 = v$  and  $u_1 = u$ ).<sup>3</sup>  
184 We can now construct a DFS traversal on  $G$  by ‘gluing’ together the sequences as,

$$(v_1, u_1, \dots, u_k, v_2, \dots, v_m) \quad (3.5)$$

185 We can run another (stochastic) DFS on  $G_1$  from  $u$  to obtain  $(u_{\tilde{\pi}(1)}, \dots, u_{\tilde{\pi}(k)})$  where  $\pi \in S_k$  and  
186  $\tilde{\pi}(1) = 1$ . We can construct a second DFS sequence as in Equation 3.5,

$$(v_1, u_{\tilde{\pi}(1)}, \dots, u_{\tilde{\pi}(k)}, v_2, \dots, v_m) \quad (3.6)$$

187 We have created two valid DFS sequences that both terminate at  $v_m$ .  $\square$

188 See Appendix B for the full details and the case of 2-edge connected graphs.

## 189 4 Related Work

190 In this section we discuss several relevant topics to graph generation. For a comprehensive review on  
191 graph generation see [28, 71].

192 **One Shot Generation** Classic generative architectures (e.g. Variational autoencoders (VAEs) [39],  
193 Generative adversarial networks (GANs) [26], etc.) work by learning a continuous mapping from  
194 a latent distribution to generate new examples with similar properties to the training distribution.  
195 These models usually incorporate a neural architecture that maps directly from the latent space to  
196 the domain of the training data (e.g. images) and therefore the output space must be predetermined.  
197 These properties pose a challenge when applied to the domain of graphs, as the latter are discrete  
198 objects with variable size and no naturally induced order. In order to circumvent these caveats, prior  
199 work [3, 17, 19, 21, 22, 29, 34, 43, 44, 55, 58, 72] has used a one-shot generation strategy. That is,  
200 the output space is limited by design to a specific representation of graphs (i.e. adjacency matrix  
201 or adjacency list) of specific size and the output is generated in a single forward pass. While the  
202 one-shot strategy has its merits, there are a few significant drawbacks such as the inability to generate  
203 graphs with arbitrary large number of nodes.

204 **Sequential Generation** The idea of using autoregressive models for graph generation is not new  
205 and there have been several works in this vein. GraphRNN [67] proposes generating BFS trajectories  
206 in order to limit the number of possible orderings per graph. Other works take a different approach  
207 of generating edges in an autoregressive manner [4, 27]. Additional approaches include Molecular-  
208 RNN [50] which incorporates a reinforcement learning environment to generate nodes and edges  
209 sequentially. Yet another approach includes sequentially generating subgraph structures [36, 41, 47].  
210 Another recent work [10] treats the induced order as a problem of dimensionality reduction and  
211 attempts to learn mappings from graphs to sequences. In this work we argue that the most effective  
212 inductive bias for the use of autoregressive models to generate graphs is to be invariant to different  
213 orderings possible under the training distribution.

214 **Molecule Generation** One of the most prominent uses of graph generation, which is used for  
215 evaluation in this work, is that of molecule generation. Molecular generation is applicable to the  
216 development of synthetic materials, drug development and more. Molecules are 3D objects which  
217 are naturally represented as point clouds<sup>4</sup> with corresponding geometric approaches [24, 56, 57, 35]  
218 which utilize inherent symmetries in the architectures employed. While 3D representations are  
219 richer and carry significant information that does not transfer to 1D and 2D representations, they

<sup>3</sup> $k$  and  $m$  denote the size of partitions and satisfy  $k + m = |\mathcal{V}|$ .

<sup>4</sup>In a point-cloud representation of a molecule each point represents an atom and bonds are implicit from the distances between atoms.

220 are costly to obtain and therefore the corresponding amount of data is limited as compared to 1D  
221 and 2D representations, which are ubiquitous. Another aspect of molecule generation is when the  
222 generation is conditioned to satisfy certain properties. For example, [59, 70, 53] generate molecules  
223 that are conditioned to bind to specific ligand structures, [38, 69] generate molecules that fulfill  
224 certain chemical properties. In this work we consider the task of de-novo generation [2, 42, 48, 61]  
225 where the objective is to generate molecules with similar properties to those in the training data.

226 **Permutation Invariant Recurrent Models** Another relevant topic is the use of autoregressive  
227 models for problems over sets which, like graphs, lack a natural order. There have been many works  
228 focusing on problems over sets. The most prominent of these is DeepSets [68] which applies a  
229 deep neural network on each element of the set and then aggregates the result with a permutation  
230 invariant operator (e.g. sum or max), finally applying another deep neural network on the aggregated  
231 result. There have also been autoregressive works designed for sets: [46] use RNNs on different  
232 permutations and output the average. While this requires  $n!$  orderings for a set of size  $n$ , the authors  
233 have presented several approximation techniques and justified them empirically. [14] have shown  
234 that while DeepSets are universal, some permutation invariant functions require unbounded width to  
235 implement successfully and have proposed using RNNs with a regularization term which enforces  
236 permutation invariance.

## 237 5 Experiments

238 A prominent application of graph generation is that of molecule design. Graph generation tasks range  
239 from de novo generation where the objective is to generate molecules with similar properties to a  
240 given dataset, to conditional generation for which the task is to generate a graph given a second graph  
241 with specific characteristics, i.e. a ligand that binds to a specific target. Our empirical evaluation  
242 focuses on the former. We evaluate our proposed regularization method on the MOSES benchmark  
243 [49] and compare to relevant baselines. Our implementation is based on the work of CharRNN which  
244 use three layers of the LSTM architecture each with hidden dimension of 600 (for complete details  
245 refer to [54]). We find a consistent improvement when adding OLR to the objective of autoregressive  
246 models.

247 The data curated by [49] is refined from the ZINC dataset [60] which contains approximately  $4.6M$   
248 molecules. The authors filter the data based on molecular weights, number of rotational bonds,  
249 lipophilicity, etc. to result in a total of  $2M$  molecules. The authors provide partitions of the data into  
250 train, test and scaffold test to allow fair evaluation.<sup>5</sup>

### 251 5.1 Computing Trajectories

252 OLR works by feeding two different trajectories that terminate at the same node. While this calculation  
253 is feasible to perform during the forward pass it introduces a computational bottleneck. In order to  
254 circumvent this issue we employ the following calculations offline. For each molecule we first index  
255 all min-cuts and randomly select one. We then generate multiple (10) traversals terminating at the  
256 same node as described in Section 3.3 and write the sequences into a file along with the original  
257 molecule from which the trajectories are derived from. When loading the data, two trajectories are  
258 selected at random and used as inputs to the OLR objective described in Section 3.2.

### 259 5.2 Data Filtering

260 Our offline computation of trajectories in Section 5.1 requires that there are min-cuts that induce  
261 sufficient number of different DFS traversals terminating at the same node. While 99.9% of the  
262 molecules in MOSES have at least two such trajectories, we filter the data to remain with molecules  
263 which have at least 10 different trajectories satisfying the criteria defined. After filtering we are left  
264 with approximately 500K molecules for training, and 55K for test and scaffold test partitions. We  
265 note that in following sections we show our method is most effective when training data is scarce and  
266 therefore the filtering process does not limit the applicability of our proposed regularization scheme.

---

<sup>5</sup>The scaffold of a molecule is the structure induced by its ring systems along with the connectivity pattern between these systems. The scaffold test partition contains molecules with structures that did not appear in the train and test partitions. The scaffold test allows for the evaluation of how well the model can generate previously unobserved scaffolds.

Table 1: Generation results at validity threshold of 0.8. Leading result highlighted in bold for each metric. Rank Average is the average position of each method over all metrics considered. As can be seen, OLR outperforms the baselines considered. Refer to the text for further details.

	Canonical	Randomized	OLR + Randomized
Unique@1K ( $\uparrow$ )	<b>1.0</b>	<b>1.0</b>	<b>1.0</b>
Unique@10K ( $\uparrow$ )	0.9965	0.9975	<b>0.9981</b>
FCD/Test ( $\downarrow$ )	<b>0.6623</b>	0.8568	0.7784
FCD/TestSF ( $\downarrow$ )	1.1980	0.4967	<b>0.4936</b>
SNN/Test ( $\uparrow$ )	0.5012	0.9947	<b>0.9958</b>
SNN/TestSF ( $\uparrow$ )	0.4835	<b>0.8246</b>	0.8220
Frag/Test ( $\uparrow$ )	0.9960	<b>1.4236</b>	1.3089
Frag/TestSF ( $\uparrow$ )	0.9943	<b>0.4795</b>	0.4769
Scaf/Test ( $\uparrow$ )	0.8193	0.9919	<b>0.9926</b>
Scaf/TestSF ( $\uparrow$ )	0.0862	<b>0.1185</b>	0.0931
IntDiv ( $\uparrow$ )	0.8515	0.8508	<b>0.8537</b>
IntDiv2 ( $\uparrow$ )	0.8457	0.8449	<b>0.8479</b>
Filters ( $\uparrow$ )	<b>0.9720</b>	0.9705	0.9702
Novelty ( $\uparrow$ )	0.9749	0.9797	<b>0.9809</b>
<b>Rank Average</b>	2.42	1.85	<b>1.50</b>

### 267 5.3 Results

268 Our results for training with OLR compared to other baselines trained on the same data are shown in  
 269 Table 1. The most relevant baseline is CharRNN [54] which is an autoregressive model trained on  
 270 Canonical SMILES. We further compare to a randomized version of CharRNN inspired by the finding  
 271 of [2] which show that augmenting the data by using randomly generated SMILES representations  
 272 of the same molecule improves performance. We also attempted to compare our method to other  
 273 non-autoregressive models such as those based on Variational Autoencoders (VAEs) [8, 25, 37];  
 274 however, we found that the models did not produce valid molecules when trained with 1000 examples,  
 275 so we do not report these results. We use the metrics defined by the MOSES benchmark [49]; see  
 276 Appendix A for a thorough description of these metrics.

277 In order to demonstrate the effectiveness of OLR we use 1000 randomly sampled data points from  
 278 the training set and evaluate over the entire test set. When training with small amounts of data there  
 279 is a trade-off between the validity of the generated molecules and the uniqueness and other metrics.  
 280 Our evaluation considers the best performing models for each method providing the validity of the  
 281 generated molecules exceeds 80%.

282 Results are depicted in Table 1. As can be seen, adding randomized variants of the molecules  
 283 outperforms the original work of [54] which train an RNN as a language model using only canonical  
 284 SMILES. Furthermore, adding the OLR objective exceeds the performance of randomized SMILES.  
 285 In order to clearly depict the performance difference, we calculate the rank of each method on each  
 286 metric considered. The average rank of each method is added as the last row of Table 1.

## 287 6 Conclusions

288 In this work we highlight the innate gap that every autoregressive model for graph generation must  
 289 mitigate - the induced order on graphs. We propose a different approach to previous works by  
 290 introducing a novel regularization scheme that encourages learning hypotheses that are invariant to  
 291 different DFS orderings. We demonstrate empirically that our proposed method improves performance  
 292 for autoregressive models and is especially effective when the available datasets are small, as is the  
 293 case in many real world problems. We believe that our approach can contribute to the applicability of  
 294 autoregressive models for graph generation and that similar ideas may be incorporated in various  
 295 generation strategies beyond the scope of this work.

## References

- 296
- 297 [1] Miltiadis Allamanis, Marc Brockschmidt, and Mahmoud Khademi. Learning to represent  
298 programs with graphs. *arXiv preprint arXiv:1711.00740*, 2017.
- 299 [2] Josep Arús-Pous, Simon Viet Johansson, Oleksii Prykhodko, Esben Jannik Bjerrum, Christian  
300 Tyrchan, Jean-Louis Reymond, Hongming Chen, and Ola Engkvist. Randomized smiles strings  
301 improve the quality of molecular generative models. *Journal of cheminformatics*, 11(1):1–13,  
302 2019.
- 303 [3] Rim Assouel, Mohamed Ahmed, Marwin H Segler, Amir Saffari, and Yoshua Bengio. De-  
304 factor: Differentiable edge factorization-based probabilistic graph generation. *arXiv preprint*  
305 *arXiv:1811.09766*, 2018.
- 306 [4] Davide Bacciu, Alessio Micheli, and Marco Podda. Edge-based sequential graph generation  
307 with recurrent neural networks. *Neurocomputing*, 416:177–189, 2020.
- 308 [5] Jesse Beisegel, Carolin Denkert, Ekkehard Köhler, Matjaz Krnc, Nevena Pivac, Robert Scheffler,  
309 and Martin Strehler. On the end-vertex problem of graph searches. *Discrete Mathematics &*  
310 *Theoretical Computer Science*, 21, 2019.
- 311 [6] Guy W Bemis and Mark A Murcko. The properties of known drugs. 1. molecular frameworks.  
312 *Journal of medicinal chemistry*, 39(15):2887–2893, 1996.
- 313 [7] Pavol Bielik, Veselin Raychev, and Martin Vechev. Phog: probabilistic model for code. In  
314 *International conference on machine learning*, pages 2933–2942. PMLR, 2016.
- 315 [8] Thomas Blaschke, Marcus Olivecrona, Ola Engkvist, Jürgen Bajorath, and Hongming Chen.  
316 Application of generative autoencoder in de novo molecular design. *Molecular informatics*,  
317 37(1-2):1700123, 2018.
- 318 [9] Marc Brockschmidt, Miltiadis Allamanis, Alexander L Gaunt, and Oleksandr Polozov. Genera-  
319 tive code modeling with graphs. *arXiv preprint arXiv:1805.08490*, 2018.
- 320 [10] Jie Bu, Kazi Sajeed Mehrab, and Anuj Karpatne. Let there be order: Rethinking ordering in  
321 autoregressive graph generation. *arXiv preprint arXiv:2305.15562*, 2023.
- 322 [11] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared  
323 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
324 language models trained on code.(2021). *arXiv preprint arXiv:2107.03374*, 2021.
- 325 [12] Xiaohui Chen, Xu Han, Jiajing Hu, Francisco JR Ruiz, and Liping Liu. Order matters: Prob-  
326 abilistic modeling of node sequence for graph generation. *arXiv preprint arXiv:2106.06189*,  
327 2021.
- 328 [13] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation  
329 of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*,  
330 2014.
- 331 [14] Edo Cohen-Karlik, Avichai Ben David, and Amir Globerson. Regularizing towards permutation  
332 invariance in recurrent models. *Advances in Neural Information Processing Systems*, 33:18364–  
333 18374, 2020.
- 334 [15] Hanjun Dai, Yingtao Tian, Bo Dai, Steven Skiena, and Le Song. Syntax-directed variational  
335 autoencoder for structured data. *arXiv preprint arXiv:1802.08786*, 2018.
- 336 [16] Alex O Davies, Nirav S Ajmeri, et al. Hierarchical gnns for large graph generation. *arXiv*  
337 *preprint arXiv:2306.11412*, 2023.
- 338 [17] Nicola De Cao and Thomas Kipf. Molgan: An implicit generative model for small molecular  
339 graphs. *arXiv preprint arXiv:1805.11973*, 2018.
- 340 [18] Jörg Degen, Christof Wegscheid-Gerlach, Andrea Zaliani, and Matthias Rarey. On the art of  
341 compiling and using ‘drug-like’ chemical fragment spaces. *ChemMedChem: Chemistry Enabling*  
342 *Drug Discovery*, 3(10):1503–1507, 2008.

- 343 [19] Yuanqi Du, Xiaojie Guo, Hengning Cao, Yanfang Ye, and Liang Zhao. Disentangled spatiotem-  
344 poral graph generative models. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
345 volume 36, pages 6541–6549, 2022.
- 346 [20] Yuanqi Du, Xiaojie Guo, Amarda Shehu, and Liang Zhao. Interpretable molecular graph gener-  
347 ation via monotonic constraints. In *Proceedings of the 2022 SIAM International Conference on*  
348 *Data Mining (SDM)*, pages 73–81. SIAM, 2022.
- 349 [21] Shuangfei Fan and Bert Huang. Labeled graph generative adversarial networks. *arXiv preprint*  
350 *arXiv:1906.03220*, 2019.
- 351 [22] Daniel Flam-Shepherd, Tony Wu, and Alan Aspuru-Guzik. Graph deconvolutional generation.  
352 *arXiv preprint arXiv:2002.07087*, 2020.
- 353 [23] Daniel Flam-Shepherd, Kevin Zhu, and Alán Aspuru-Guzik. Language models can learn  
354 complex molecular distributions. *Nature Communications*, 13(1):3293, 2022.
- 355 [24] Victor Garcia Satorras, Emiel Hooeboom, Fabian Fuchs, Ingmar Posner, and Max Welling.  
356 E (n) equivariant normalizing flows. *Advances in Neural Information Processing Systems*,  
357 34:4181–4192, 2021.
- 358 [25] Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato,  
359 Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel,  
360 Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven  
361 continuous representation of molecules. *ACS central science*, 4(2):268–276, 2018.
- 362 [26] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
363 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications*  
364 *of the ACM*, 63(11):139–144, 2020.
- 365 [27] Nikhil Goyal, Harsh Vardhan Jain, and Sayan Ranu. Graphgen: a scalable approach to domain-  
366 agnostic labeled graph generation. In *Proceedings of The Web Conference 2020*, pages 1253–  
367 1263, 2020.
- 368 [28] Xiaojie Guo and Liang Zhao. A systematic survey on deep generative models for graph  
369 generation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- 370 [29] Xiaojie Guo, Liang Zhao, Zhao Qin, Lingfei Wu, Amarda Shehu, and Yanfang Ye. Interpretable  
371 deep graph generation with node-edge co-disentanglement. In *Proceedings of the 26th ACM*  
372 *SIGKDD international conference on knowledge discovery & data mining*, pages 1697–1707,  
373 2020.
- 374 [30] Vincent J Hellendoorn, Charles Sutton, Rishabh Singh, Petros Maniatis, and David Bieber.  
375 Global relational models of source code. In *International conference on learning representations*,  
376 2019.
- 377 [31] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
378 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in*  
379 *neural information processing systems*, 30, 2017.
- 380 [32] Abram Hindle, Earl T Barr, Mark Gabel, Zhendong Su, and Premkumar Devanbu. On the  
381 naturalness of software. *Communications of the ACM*, 59(5):122–131, 2016.
- 382 [33] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*,  
383 9(8):1735–1780, 1997.
- 384 [34] Shion Honda, Hirotaka Akita, Katsuhiko Ishiguro, Toshiki Nakanishi, and Kenta Oono. Graph  
385 residual flow for molecular graph generation. *arXiv preprint arXiv:1909.13521*, 2019.
- 386 [35] Emiel Hooeboom, Victor Garcia Satorras, Clément Vignac, and Max Welling. Equivariant  
387 diffusion for molecule generation in 3d. In *International conference on machine learning*, pages  
388 8867–8887. PMLR, 2022.

- 389 [36] Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Junction tree variational autoencoder for  
390 molecular graph generation. In *International conference on machine learning*, pages 2323–2332.  
391 PMLR, 2018.
- 392 [37] Artur Kadurin, Alexander Aliper, Andrey Kazennov, Polina Mamoshina, Quentin Vanhaelen,  
393 Kuzma Khrabrov, and Alex Zhavoronkov. The cornucopia of meaningful leads: Applying deep  
394 adversarial autoencoders for new molecule development in oncology. *Oncotarget*, 8(7):10883,  
395 2017.
- 396 [38] Seokho Kang and Kyunghyun Cho. Conditional molecular design with deep generative models.  
397 *Journal of chemical information and modeling*, 59(1):43–52, 2018.
- 398 [39] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*  
399 *arXiv:1312.6114*, 2013.
- 400 [40] Greg Landrum. Rdkit: Open-source cheminformatics. 2006. *Google Scholar*, 2006.
- 401 [41] Renjie Liao, Yujia Li, Yang Song, Shenlong Wang, Will Hamilton, David K Duvenaud, Raquel  
402 Urtasun, and Richard Zemel. Efficient graph generation with graph recurrent attention networks.  
403 *Advances in neural information processing systems*, 32, 2019.
- 404 [42] Jaechang Lim, Seongok Ryu, Jin Woo Kim, and Woo Youn Kim. Molecular generative  
405 model based on conditional variational autoencoder for de novo molecular design. *Journal of*  
406 *cheminformatics*, 10(1):1–9, 2018.
- 407 [43] Tengfei Ma, Jie Chen, and Cao Xiao. Constrained generation of semantically valid graphs via  
408 regularizing variational autoencoders. *Advances in Neural Information Processing Systems*, 31,  
409 2018.
- 410 [44] Kaushalya Madhawa, Katushiko Ishiguro, Kosuke Nakago, and Motoki Abe. Graphnvp: An  
411 invertible flow model for generating molecular graphs. *arXiv preprint arXiv:1905.11600*, 2019.
- 412 [45] Andreas Mayr, Günter Klambauer, Thomas Unterthiner, Marvin Steijaert, Jörg K Wegner, Hugo  
413 Ceulemans, Djork-Arné Clevert, and Sepp Hochreiter. Large-scale comparison of machine  
414 learning methods for drug target prediction on chembl. *Chemical science*, 9(24):5441–5451,  
415 2018.
- 416 [46] Ryan L Murphy, Balasubramaniam Srinivasan, Vinayak Rao, and Bruno Ribeiro. Janossy  
417 pooling: Learning deep permutation-invariant functions for variable-size inputs. *arXiv preprint*  
418 *arXiv:1811.01900*, 2018.
- 419 [47] Marco Podda, Davide Bacciu, and Alessio Micheli. A deep generative model for fragment-based  
420 molecule generation. In *International Conference on Artificial Intelligence and Statistics*, pages  
421 2240–2250. PMLR, 2020.
- 422 [48] Peter Pogány, Navot Arad, Sam Genway, and Stephen D Pickett. De novo molecule design  
423 by translating from reduced graphs to smiles. *Journal of chemical information and modeling*,  
424 59(3):1136–1146, 2018.
- 425 [49] Daniil Polykovskiy, Alexander Zhebrak, Benjamin Sanchez-Lengeling, Sergey Golovanov,  
426 Oktai Tatanov, Stanislav Belyaev, Rauf Kurbanov, Aleksey Artamonov, Vladimir Aladinskiy,  
427 Mark Veselov, et al. Molecular sets (moses): a benchmarking platform for molecular generation  
428 models. *Frontiers in pharmacology*, 11:565644, 2020.
- 429 [50] Mariya Popova, Mykhailo Shvets, Junier Oliva, and Olexandr Isayev. Molecularrrn: Generating  
430 realistic molecular graphs with optimized properties. *arXiv preprint arXiv:1905.13372*, 2019.
- 431 [51] Oleksii Prykhodko, Simon Viet Johansson, Panagiotis-Christos Kotsias, Josep Arús-Pous,  
432 Esben Jannik Bjerrum, Ola Engkvist, and Hongming Chen. A de novo molecular generation  
433 method using latent vector based generative adversarial network. *Journal of Cheminformatics*,  
434 11(1):1–13, 2019.
- 435 [52] David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of chemical*  
436 *information and modeling*, 50(5):742–754, 2010.

- 437 [53] Eyal Rozenberg and Daniel Freedman. Semi-equivariant conditional normalizing flows, with  
438 applications to target-aware molecule generation. *Machine Learning: Science and Technology*,  
439 4(3):035037, 2023.
- 440 [54] Marwin HS Segler, Thierry Kogej, Christian Tyrchan, and Mark P Waller. Generating focused  
441 molecule libraries for drug discovery with recurrent neural networks. *ACS central science*,  
442 4(1):120–131, 2018.
- 443 [55] Chence Shi, Minkai Xu, Zhaocheng Zhu, Weinan Zhang, Ming Zhang, and Jian Tang.  
444 Graphaf: a flow-based autoregressive model for molecular graph generation. *arXiv preprint*  
445 *arXiv:2001.09382*, 2020.
- 446 [56] Gregor Simm, Robert Pinsler, and José Miguel Hernández-Lobato. Reinforcement learning  
447 for molecular design guided by quantum mechanics. In *International Conference on Machine*  
448 *Learning*, pages 8959–8969. PMLR, 2020.
- 449 [57] Gregor NC Simm, Robert Pinsler, Gábor Csányi, and José Miguel Hernández-Lobato.  
450 Symmetry-aware actor-critic for 3d molecular design. *arXiv preprint arXiv:2011.12747*, 2020.
- 451 [58] Martin Simonovsky and Nikos Komodakis. Graphvae: Towards generation of small graphs  
452 using variational autoencoders. In *Artificial Neural Networks and Machine Learning–ICANN*  
453 *2018: 27th International Conference on Artificial Neural Networks, Rhodes, Greece, October*  
454 *4-7, 2018, Proceedings, Part I 27*, pages 412–422. Springer, 2018.
- 455 [59] Miha Skalic, Davide Sabbadin, Boris Sattarov, Simone Sciabola, and Gianni De Fabritiis. From  
456 target to drug: generative modeling for the multimodal structure-based ligand design. *Molecular*  
457 *pharmaceutics*, 16(10):4282–4291, 2019.
- 458 [60] Teague Sterling and John J Irwin. Zinc 15–ligand discovery for everyone. *Journal of chemical*  
459 *information and modeling*, 55(11):2324–2337, 2015.
- 460 [61] Xiaochu Tong, Xiaohong Liu, Xiaoqin Tan, Xutong Li, Jiabin Jiang, Zhaoping Xiong, Tingyang  
461 Xu, Hualiang Jiang, Nan Qiao, and Mingyue Zheng. Generative models for de novo drug design.  
462 *Journal of Medicinal Chemistry*, 64(19):14011–14027, 2021.
- 463 [62] Leonid Nisonovich Vaserstein. Markov processes over denumerable products of spaces, de-  
464 scribing large systems of automata. *Problemy Peredachi Informatsii*, 5(3):64–72, 1969.
- 465 [63] Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. Order matters: Sequence to sequence for  
466 sets. *arXiv preprint arXiv:1511.06391*, 2015.
- 467 [64] David Weininger. Smiles, a chemical language and information system. 1. introduction to  
468 methodology and encoding rules. *Journal of chemical information and computer sciences*,  
469 28(1):31–36, 1988.
- 470 [65] Pengcheng Yin and Graham Neubig. A syntactic neural model for general-purpose code  
471 generation. *arXiv preprint arXiv:1704.01696*, 2017.
- 472 [66] Jiaxuan You, Bowen Liu, Zhitao Ying, Vijay Pande, and Jure Leskovec. Graph convolutional  
473 policy network for goal-directed molecular graph generation. *Advances in neural information*  
474 *processing systems*, 31, 2018.
- 475 [67] Jiaxuan You, Rex Ying, Xiang Ren, William Hamilton, and Jure Leskovec. Graphrnn: Generat-  
476 ing realistic graphs with deep auto-regressive models. In *International conference on machine*  
477 *learning*, pages 5708–5717. PMLR, 2018.
- 478 [68] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov,  
479 and Alexander J Smola. Deep sets. *Advances in neural information processing systems*, 30,  
480 2017.
- 481 [69] Chengxi Zang and Fei Wang. Moflow: an invertible flow model for generating molecular graphs.  
482 In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery &*  
483 *data mining*, pages 617–626, 2020.

- 484 [70] Zaixi Zhang, Yaosen Min, Shuxin Zheng, and Qi Liu. Molecule generation for target pro-  
 485 tein binding with structural motifs. In *The Eleventh International Conference on Learning*  
 486 *Representations*, 2022.
- 487 [71] Yanqiao Zhu, Yuanqi Du, Yinkai Wang, Yichen Xu, Jieyu Zhang, Qiang Liu, and Shu Wu. A  
 488 survey on deep graph generation: Methods and applications. *arXiv preprint arXiv:2203.06714*,  
 489 2022.
- 490 [72] Dongmian Zou and Gilad Lerman. Encoding robust representation for graph generation. In  
 491 *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–9. IEEE, 2019.

## 492 A Metric Details

493 In this section we provide details for the metrics reported in Table 1.

494 A few of the similarity measures (SNN and IntDiv) are based on the *Tanimoto coefficient*. In order to  
 495 compute the Tanimoto coefficient, the molecules are mapped to a vector of fingerprints where each  
 496 bit in the vector represents the presence (or absence) of a specific fragment.<sup>6</sup> For molecules  $A, B$ ,  
 497 denote their fingerprints by  $m_A$  and  $m_B$  respectively, the Tanimoto coefficient is then calculated as  
 498 the Jaccard index of the two vectors,

$$J(m_A, m_B) = \frac{|m_A \cap m_B|}{|m_A \cup m_B|} = \frac{|m_A \cap m_B|}{|m_A| + |m_B| - |m_A \cap m_B|}. \quad (\text{A.1})$$

499 We denote the Tanimoto coefficient of molecules  $A, B$  by  $T(A, B)$ .

500 **Unique@K** report the fraction of uniquely generated valid SMILES strings amongst the  $K$   
 501 molecules generated (validity is determined by the RDKit library). We generate 30,000 molecules  
 502 for each model and report for  $K = 1,000$  and  $K = 10,000$ . High uniqueness values ensure the  
 503 models do not collapse into repeatedly producing the same set of molecules.

504 **Fréchet ChemNet Distance (FCD)** is a metric for evaluating generative models in the chemical  
 505 context, the method is based on the well established *Fréchet Inception Distance* (FID) metric used to  
 506 evaluate the performance of generative models in computer vision [31].

507 Fréchet distance measure the Wasserstein-2 distance [62] from the distributions induced by taking  
 508 the activations of the last layer of a relevant deep neural net. In the case of FCD, molecule activations  
 509 are probed from ChemNet [45]. Given a set of generated molecules, denote by  $G$  the set of vectors  
 510 as obtained by the activations of ChemNet, one can calculate the mean and covariance  $\mu_G$  and  $\Sigma_G$ .  
 511 Similarly, denote  $\mu_R$  and  $\Sigma_R$  the mean and covariance of the set of molecules in the reference set,  
 512 the FCD is calculated as follows,

$$FCD(G, R) = \|\mu_G - \mu_R\|^2 + Tr\left(\Sigma_G + \Sigma_R - 2(\Sigma_G \Sigma_R)^{1/2}\right). \quad (\text{A.2})$$

513 where  $Tr(M)$  denotes the trace of the matrix  $M$ . Low FCD values indicate that the generated  
 514 molecules distribute similarly to the reference set.

515 **Similarity to Nearest Neighbor (SNN)** is the average of the Tanimoto coefficient of the generated  
 516 molecule set denoted by  $G$  and their respective nearest neighbor in a reference set of molecules  
 517 denote by  $R$ . High SNN indicates the generated molecules have similar structures to those in the  
 518 reference set. This metric is in the range of  $[0, 1]$ .

519 **Fragment similarity (Frag)** is a fragment similarity measure based on the BRICS fragments [18].  
 520 Denote the set of BRICS fingerprints vectors of the generated molecules by  $G$  and similarly  $R$  for the  
 521 reference molecules. The fragment similarity is defined as the cosine similarity of the sum vectors,

$$Frag(G, R) = cosine\left(\sum_{g \in G} g, \sum_{r \in R} r\right) \quad (\text{A.3})$$

<sup>6</sup>The molecular fingerprints are obtained from RDKit [40] and are based on the extended-connectivity fingerprints [52].

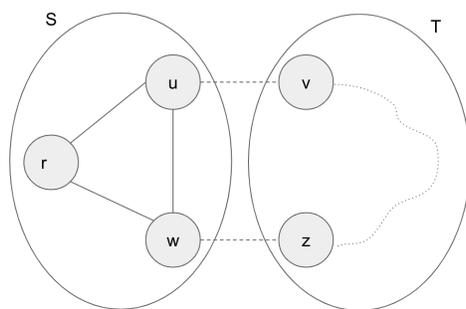


Figure 3: Proof illustration -  $S$  has a cycle and two different trajectories starting from  $u$  and ending with  $w$  ( $urw$  and  $uw(r)$ ). Concatenating with the trajectory from  $z$  to  $v$  we obtain two different DFS trajectories with a shared suffix.

522 The Frag measure is in the range of  $[0, 1]$ , values closer to 1 indicate that the generated and reference  
 523 molecule set have a similar distribution of BRICS fragment.

524 **Scaffold similarity (Scaff)** is similar to the fragment similarity, instead of the BRICS fragment,  
 525 Scaff is based on mapping molecules to their Bemis–Murcko scaffolds [6].<sup>7</sup> The measure also has a  
 526 range of  $[0, 1]$ , values closer to 1 indicate that the generated molecule set has a similar distribution of  
 527 scaffold to the reference set.

528 **Internal diversity (IntDiv)** is a measure of the chemical diversity within a generated set of  
 529 molecules  $G$ . This metric indicates

$$IntDiv_p = 1 - \left( \frac{1}{|G|^2} \sum_{A, B \in G} T(A, B)^p \right)^{1/p} \quad (\text{A.4})$$

530 We report the internal diversity for  $p = 1, 2$ . This measure is in the range  $[0, 1]$ . Low values indicate  
 531 a lack of diversity in the generated molecules, i.e. that the model outputs molecules with similar  
 532 fingerprints.

533 **Filters** is the fraction of generated molecules that pass a certain filtering that has been applied to  
 534 the training data. The metric is in the range of  $[0, 1]$ , high values indicate that the model has learnt to  
 535 generate molecules which avoid the structures omitted by the filtering process.

536 **Novelty** is the fraction of generated molecules that does not appear in the training set. This measure  
 537 is in the range of  $[0, 1]$  and is an indication of the whether the model overfits the training data.

## 538 B Missing Proofs

539 In this section we show how to construct distinct DFS trajectories with common end vertex for a  
 540 2–edge connected graph conditioned that the graph is not a cycle.

541 *Proof.* From our assumption that the graph is not a cycle, there exists at least two nodes with degree  
 542  $\geq 3$ . Denote by  $C = (S, T)$  a minimal cut of size 2 (such a cut exists from our assumption that  
 543 the graph is 2-connected). Denote the edges of the minimal cut by  $e_1 = (u, v)$  and  $e_2 = (w, z)$   
 544 such that  $u, w \in S$  and  $v, z \in T$ . Next, we claim that at least one of the partitions contains a cycle,  
 545 otherwise there is a path connecting  $S$  and  $T$  since there are nodes in the graph which have a degree  
 546 of 3 in the original graph with a path between them. Assume with out loss of generality that  $S$  is the  
 547 partition with a cycle, therefore there are at least 2 different traversals of  $S$  that start with  $u$  and end

<sup>7</sup>Bemis–Murcko scaffold is the ring structure of a molecule along with the bonds connecting the rings, i.e. the molecule without the side chains.

548 with  $w$ . There is also a trajectory between  $z$  and  $v$ . Putting together, there are at least 2 trajectories  
549 of the entire graph with a common suffix which is the traversal of  $T$ . Figure 3 illustrates the proof  
550 concept. □