# Fragment-based t-SMILES for de novo molecular generation

Juan-Ni Wu, Tong Wang, Yue Chen, Li-Juan Tang, Hai-Long Wu*, Ru-Qin Yu*

*State Key Laboratory of Chemo/Biosensing and Chemometrics, College of Chemistry and Chemical Engineering, Hunan University, Changsha 410082, People's Republic of China*

## ABSTRACT

At present, sequence-based and graph-based models are two of popular used molecular generative models. In this study, we introduce a general-purposed, fragment-based, hierarchical molecular representation named t-SMILES (tree-based SMILES) which describes molecules using a SMILES-type string obtained by doing breadth first search (BFS) on full binary molecular tree formed from fragmented molecular graph. The proposed t-SMILES combines the advantages of graph model paying more attention to molecular topology structure and language model possessing powerful learning ability. Experiments with feature tree rooted JTVAE and chemical reaction-based BRICS molecular decomposing algorithms using sequence-based autoregressive generation models on three popular molecule datasets including Zinc, QM9 and ChEMBL datasets indicate that t-SMILES based models significantly outperform previously proposed fragment-based models and being competitive with classical SMILES based and graph-based approaches. Most importantly, we proposed a new perspective for fragment based molecular designing. Hence, SOTA powerful sequence-based solutions could be easily applied for fragment based molecular tasks.

**Keywords**: Fragment-based drug discovery, Tree-based SMILES, Sequence-based De novo design

## 1 Introduction

The starting point for molecular drug discovery programs is to identify initial ''hit'' compounds that bind to the target and have the potential for optimization of clinical candidates with the desired therapeutic effect. Fragment-based drug discovery (FBDD) builds drugs from small molecular pieces. Since the pioneering work by Fesik and coworkers in 1996[1], FBDD has become a

recognized technology in the pharmaceutical industry and within academia. More than 30 drug candidates derived from fragments have been reported to enter the clinic[2]. Compared with atom-based techniques, the size of the search space is greatly reduced by the use of the fragment strategy. In addition, fragments could provide fundamental insights into molecular recognition between proteins and ligands. As a consequence, there is a higher probability of finding molecules that match the known targets.

The early FBDD was generally implemented based on fragment library using virtual screening and other methods. Although deep learning[3] has been widely used in molecular generation tasks[4,5] with fragment-based method as a research topic, the method of fragmenting molecules and coding molecular substructures in the form of a string-type sequence like SMILES[6] to finally generate molecules has not yet been fully explored.

In recent years, various deep generative models for the task of automatically generate molecules have been proposed. Among deep learning-based methods, models with sequence representations [7-9] such as SMILES and 2D representations such as graphs [10-14] are most popular, while recently a plethora of models generating molecules in 3D[15] also starts to attract attention.

As a more natural representation of molecules, generally speaking, graph neural network could generate 100% valid molecules as it can easily implement valence bond constraints and other verification rules. However, it has been shown that the expressive power of standard GNNs is bounded by Weisfeiler-Leman (WL) graph isomorphism phenomenon, the lack of ways to model long-range interactions and higher-order structures limited the use of GNNs[16], though some recent studies have proposed new methods such as subgraph isomorphism[17], message-passing simple networks[18] and many others techniques to improve the expressive power of standard GNNs[19].

From the perspective of graph-based computational procedure, SMILES is a linear string obtained by performing depth first search (DFS) on molecular graph, which is more like human natural language. When generating SMILES, the chemical graph is firstly trimmed to remove hydrogen atoms and cycles are broken turning them into a spanning tree. Where cycles have been broken, numeric suffix labels are included to indicate the connected nodes. Parentheses are used to indicate points of branching on the tree. The generation algorithm of classical SMILES directly breaks down the most common ring structures in molecules. As a consequence, some elements of SMILES syntax must occur in pairs with deep nesting to represent molecular topological structure. Without discussing what chemical information could be learned, models trained on SMILES somehow generate part chemical invalid strings, particularly when trained on small datasets, which some have identified as a limitation need to be addressed[20]. Two alternative solutions to the classical SMILES have been proposed[20]. The DeepSMILES[21] aims to remove long term dependencies associated with the representation of rings and branches from the SMILES syntax to finally increase the proportion of valid molecules generated. Self-referencing embedded strings (SELFIES)[22] are an entirely different molecular representation based on a Chomsky type-2 grammar, in which every SELFIES string specifies a valid chemical graph.

In addition, almost all substructure (motif or fragment) based methods published to date rely on a substructure dictionary (motif vocabulary dictionary or fragment library) of candidate fragments[13,23-28]. In molecules, a small group of fragments being used frequently, while most of fragments are rarely used. The differences among dictionary-based approaches arise solely from how the database is searched, or the contents of the database itself. As a result, these methods are inherently constrained to a set of predetermined rules or examples, limiting the exploration of chemical space and the learning ability of the models to a certain extent.

Recently, attention-based Transformer[29] pre-trained models have been proved to enable text generation with human-like capabilities, including texts with specific properties such as style or subject, if trained on enough data. With the rapid development of natural language processing(NLP) technology, and the increasing interest in larger and more complex molecules for treatment, language models may show a better ability to learn complex molecules than most graph generation models[30].

Motivated by the success of NLP and the strategies of FBDD, we hope to adopt sequence-based models to handle fragment-based molecular generation tasks. So that, we propose a new molecule code t-SMILES based on fragmented molecule, which describes a molecule with classical SMILES-type string and takes language model as the main generation model. t-SMILES based model combines the advantages of graph model paying attention to molecular topology structure and language model possessing powerful NLP learning ability.

The t-SMILES method firstly generates an acyclic molecular tree (AMT) whose role is to represent fragmented molecules. Fragments in this tree are chemically valid substructures automatically extracted from molecules using a molecular disconnection algorithm. In the second stage, the AMT is transformed to a full binary tree (FBT). Finally, breadth-first traversal of the FBT yields a t-SMILES string.

Our proposed method is a general framework that does not limit the substructure scheme as long as it could be used to generate chemically valid fragments and construct valid AMT. t-SMILES string could be directly applied to sequence-based machine learning models without adjusting the model architecture in specific fields. In addition, due to the hierarchical representation of molecules, our model can clearly learn the high-level topology structural information of molecules while processing the detailed substructure information. Most importantly, by exploiting NLP

4

methods, t-SMILES based models opens an exploration door for fragment-based molecule tasks. Therefore, SOTA sequence-based language models could be easily used for fragment-based molecular tasks.

To evaluate the feasibility and adaptability of t-SMILES, we break down molecules using two different strategies BRICS[31] based and feature tree rooted JTVAE[12] based algorithms respectively, and then use Transformer-decoder based autoregressive models to generate molecules on three popular molecular datasets Zinc[32], QM9[33] and ChEMBL[34]. Finally, we compare t-SMILES based models with the most popular graph neural network models, SMILES-based models and other baselines.

Five metrics including validity, uniqueness, novelty, KL divergence(KLD)[35]and Fréchet ChemNet distance score(FCD)[36] introduced in GuacaMol[37] benchmark are used to evaluate the general performance of the models. In addition, we use logP, Penalized logP(plogP)[38], SA score(SAS)[39], BertzCT[40], QED[41], TPSA, NP Score(NPS)[42], and FractionCSP3 to evaluate whether the generation model could effectively learn the physical and chemical properties of molecules.

The t-SMILES construction philosophy and detailed comparative experiments show that the validity of generated molecules by t-SMILES models could be greater than 99%. At the same time, using 5 layers mini GPT2[43](mGPT2) model, t-SMILES could well capture the physicochemical properties of molecules to maintain the similarity of the generated molecules to the distribution of dataset molecules, which make t-SMILES based models significantly outperforming previously proposed fragment-based models and being competitive with classical SMILES and graph-based approaches. In addition, compared with classical SMILES which is relatively difficult to be

augmented[20], t-SMILES is easily to be expanded to explore different chemical spaces by using different molecular fragmentation algorithms.

Although current de novo molecular generation approaches have made impressive advances, the sequence-based perspectives create new opportunities to advance fragment-based molecular design tasks. The t-SMILES solution demonstrates that de novo generation of molecules from fragmented SMILES is possible. This solution challenges the current research paradigm used for FBDD.

## 2 Methods

In order to support t-SMILES algorithm, on the basis of the classic SMILES, we introduce a new character '&' to represent the end of the sub-branch, and another new character '^' to separate the two adjacent SMILES substructure segments. In this section, we firstly introduce the general idea of t-SMILES, and then introduce FBT and AMT which are the core parts of t-SMILES algorithm, followed with molecular fragmentation algorithms and finally discuss the molecular reconstruction strategy.

### 2.1 t-SMILES Algorithm Overview

In t-SMILES algorithm, molecular graph is firstly divided into chemical valid fragments (or substructures, clusters, subgroups, subgraphs) using a specified or more disconnection methods to obtain its AMT shown in the middle of Figure 1. Following with the AMT being transformed into a FBT shown in the right of Figure 1, and finally the FBT is traversed in breadth first search (BFS) to obtain the t-SMILES string. During the reconstruction, the reverse process is used, and finally the molecular fragments are assembled into the chemical correct molecular graph.
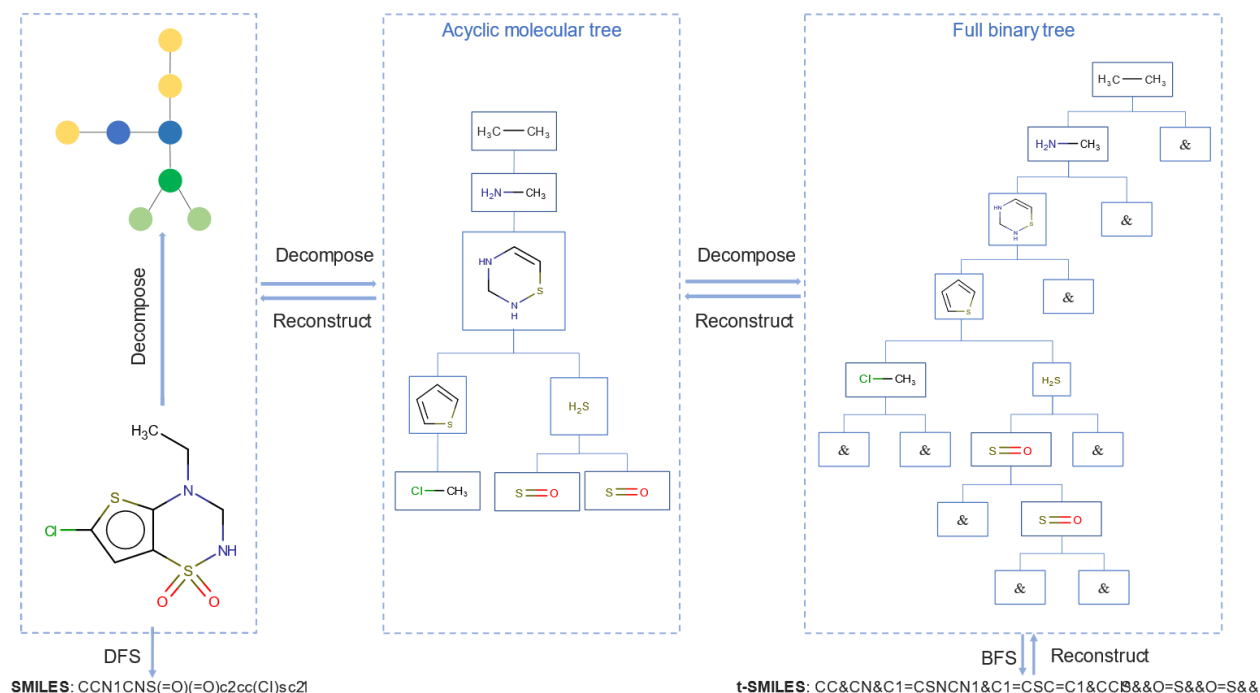
6

**Figure 1**. Overview of t-SMILES algorithm: A molecular graph G is first decomposed into its reduced graph, where each colored node in the tree represents a substructure in the molecule. We then generate an AMT based on reduced graph, following with trasnformation of AMT to FBT. Finally, the FBT is traversed in BFS to obtain its t-SMILES code. To reconstruct the molecule, we first rebuild FBT from t-SMILES string, and then transform FBT to AMT, finally assemble nodes in the tree back to the original molecule.

We follow below steps to build t-SMILES:

**Algorithm steps to construct t-SMILES from molecule**

Step 1: Break down molecule according to the selected molecular fragmentation algorith to build AMT;

Step 2: Convert the AMT to a FBT through algorithm;

Step 3: Traverse the FBT with BFS algorithm to get t-SMILES .

BFS algorithm for the tree is a level order traversal of tree. For any node w in the BFS tree rooted at v, the tree path from v to w corresponds to a shortest path from v to w in the corresponding graph.

In the following, the terms nodes and subtrees are used for describing ATM and FBT, the terms fragment, substructure and subgroup are used to describe a part of the molecule. Subtrees,

fragments, substructures or subgroups are always assumed to be connected parts. FBT has two children named left subtree and right subtree.

## 2.2 Full Binary Tree

Tree is the core concept in proposed t-SMILES algorithm. A tree is a special type of graph in which there is just a single path connecting each pair of vertices, that is, there are no cycles or rings within the graph. The root node of a tree is the starting point while the other vertices are either branch nodes or leaf nodes.

A FBT is a special type of binary tree in which every parent node/internal node has either two or no children. As the most trivial tree, FBT's structure is regular and easy to calculate. The reason for using FBT with some redundant nodes instead of complete binary tree or other trees is that its algorithm and structure being easy to learn by deep learning models, and the redundant nodes could be used as global marker nodes. In this work, the character '&'(tree node marked as '&') marking the end of the tree branch in the FBT could be regarded as the global structural information describing the molecular topology in t-SMILES string.

With chemical meaningful molecular fragmentation representation using FBT, t-SMILES effectively reduces the nesting depth of brackets in classical SMILES codes, weakens the long-term dependencies in sequences, and fundamentally reduces the difficulty of learning molecular information for using sequence-based deep learning models. In t-SMILES algorithm, except for the extra two characters '**&**' and '**^**', no more symbols are introduced, nor are recursion and other sophisticated calculations with high computational complexity introduced.

## 2.3 Molecule Decomposition

The key and first step in t-SMILES is to decompose molecules into chemical valid fragments according to a specified disconnection algorithm. According to Lounkine *et al*.[44], there exists four major strategies to fragment designing: knowledge-based, synthetically oriented, random, and systematic or hierarchical. The open-source molecular toolkit RDKit[45] has implemented some molecular fragmentation methods, such as RECAP[46] and BRICS[31] etc.

RECAP and BRICS both disconnect a molecule to fragments on bonds based on chemical reaction rules. In the RECAP method, molecules are cleaved at 11 chemical bond types that correspond to common chemical reactions, while BRICS attempts to improve RECAP for molecule fragmentation by using a more elaborate set of 16 rules and additional pre- and postfilters. In this study, considering that the training datasets are mainly small molecules, and the molecular fragments segmented by BRICS are large or molecules could not be broken down, say on the subset of Zinc and ChEMBL, if necessary, we could further cut off the branch structures connected to ring structures on the basis of BRICS, and then cut off the bridge bonds between two rings. BRICS algorithm is not used on QM9 in this study as molecules are relatively rather small to be broken down. We generate reduced graph according to the cutting off logic of BRICS and then calculate its spanning tree as AMT.

Another molecular decomposition algorithm in our study could find its root in feature tree[47] published for molecular similarity algorithm by Rarey *et al*. in 1998 and JTVAE[12] proposed later by Jin *et al.*  Given a graph G, we first find all its simple cycles, and its edges not belonging to any cycles. Two simple rings are merged together if they have more than two overlapping atoms, as they constitute a specific structure called bridged compounds [48]. Each of those cycles or edges

is considered as a cluster. Next, a cluster graph is constructed by adding edges between all intersecting clusters. Finally, we select one of its spanning trees as the AMT of G.

The t-SMILES algorithm is a general and fragment based molecular representation framework which does not limit the choice of molecular decomposing algorithms. Different molecular segmentation algorithm may require different fragment assembling algorithm to complete the reconstruction of the generated molecules.

## 2.4 Acyclic Molecular Tree

The idea of using tree as the base data structure of algorithms to address molecular related tasks has been long established in cheminformatics. In early study of molecular descriptor and similarity analysis, algorithms such as reduced graph[49], feature tree[47,50] not only had shown potential power to improve the similarity search but also being capable of retrieving more diverse active compounds than using Daylight fingerprints[51].Some recent works[52-55] proposed to incorporate tree-based deep learning models into molecular generation and synthesis tasks as well.

AMT being capable to describe the molecule at various levels of resolution, reduced graph[56] provides summary representations of chemical structures by collapsing groups of connected atoms into single nodes while preserving the topology of the original structures. In reduced graph, the nodes represent groups of atoms and the edges represent the physical bonds between the groups. Constructing reduced graph in this way forms a hierarchical graph, whose top layer being the molecular topology representing global information, and the bottom layer representing molecular fragments of detailed information. Groups of atoms are clustered into a node in the reduced graph approach, which could be done based on fragmentation algorithms. The feature tree[47] is a representation of a molecule similar to a reduced graph. The vertices of the feature tree are molecular fragments and edges connect vertices that represent fragments connected in the simple

molecular graph. In t-SMILES algorithm, the minimum spanning tree of the reduced graph and the concept of feature tree could be regarded as an AMT in the intermediate step, and then the next encoding algorithm is done based on this AMT.

Specific to our experiments, one approach is to generate AMT based on the tree logic fragmented by the BRICS algorithm. Another method uses junction tree[57] introduced by Jin *et al.* in JTVAE as the AMT. In the case of a junction tree, each node in the tree corresponds to a subset (subgraph, group, cluster or clique) in the original graph and edge is used to connect two clusters. In this study, nodes of Junction tree represent rings, bonds, bridged compounds, or singletons in the original molecular graph which are generated by fragment decomposing algorithm and edges represent the physical bonds between groups.

## 2.5 Molecular Reconstruction and Optimization

In the process of reconstructing the generated molecules from t-SMILES strings, we follow below steps:

**Algorithm steps to reconstruct molecule from t-SMILES**

Step 1: Decompose t-SMILES  to reconstruct the FBT;
Step 2: Convert the FBT to the AMT;
Step 3: Assemble molecular fragments according to the selected algorithm to generate the correct molecular graph and then optimize it.

During reconstruction, one key problem is how to assemble the molecular fragments together to get a 'chemical correct' molecule. Ideally, the assemble algorithm should be selected to match the molecular fragmentation algorithm. In this study, for efficiency reasons, we assemble molecular graph one neighborhood at a time, following the order in which the tree itself was generated. In other words, we start from the root node of AMT and its neighbors, then we proceed to assemble

the neighbors and their associated clusters, and so on. If there is more than one candidate when assembling two pieces, we simply select one randomly.

The overall performance of t-SMILES based generative model is controlled by two main factors. The first one is whether t-SMILES can be learned and generated efficiently, and the other one is the reconstruction and molecular optimization algorithm. We evaluate the assembling algorithms by directly reconstructing molecules from training set data. Detailed metrics are shown in Table 1 and Figure 2.

**Table 1.** Distributional results by directly reconstructing molecules from the training dataset.

| Dataset | Model | Valid(↑) | Uniq(↑) | Novel(↑) | KLD(↑) | FCD(↑) |
|---------|-------|----------|---------|----------|--------|--------|
| **Zinc** | JTVAE | 1.000 | 0.770 | 0.662 | 0.982 | 0.811 |
| | BRICS | 1.000 | 0.780 | 0.681 | 0.986 | 0.849 |
| **QM9** | JTVAE | 1.000 | 0.741 | 0.245 | 0.976 | 0.968 |
| **ChEMBL** | JTVAE | 1.000 | 0.786 | 0.677 | 0.969 | 0.694 |
| | BRICS | 1.000 | 0.787 | 0.672 | 0.977 | 0.792 |

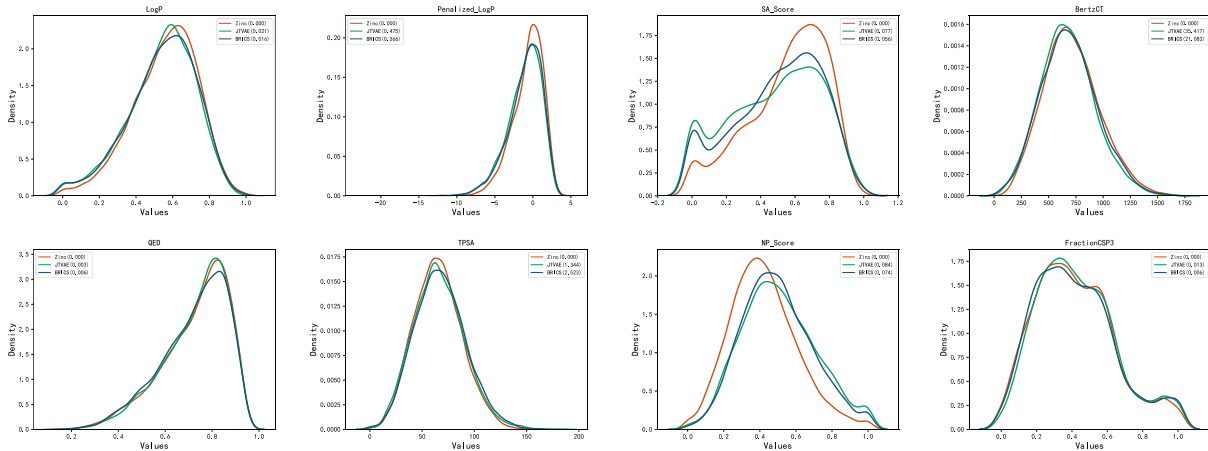Data in Table 1 show that FCD scores on ChEMBL are the lowest ones while it is the highest on QM9.



**Figure 2.** Distributions of physiochemical properties (logP, plogP, SAS, BertzCT, QED, TPSA, NPS and FractionCSP3) of reconstructed molecules directly from Zinc training set. Wasserstein distance is used in this figure and later.

From Figure 2, it can be seen that SAS and NPS are the two most variable metrics when reconstituting molecules.

If we output all possible assembly results, we could get a set of molecules which come from the same fragments group with different structures. From this point of view, generating new molecules with the desired properties (desired structure) rather than duplicating the training set is exactly the potential goal of the molecule generation task and not a negative aspect. MOG[58] argued that the common pitfall of existing molecule generation models based on distributional learning is that the exploration is confined to the training distribution, and the generated molecules exhibit "the striking similarity" with known molecules included in the training set. Models that do not require training molecules are free from this problem, but they introduce other problems such as long training time, the sensitivity of balance between exploration and exploitation, large variance, and importantly, a lack of information about the known distribution. Based on t-SMILES, it's possible to select appropriate optimization algorithm to control how the fragments are assembled, thereby controlling the properties of the output molecules. Molecule optimization is another challenging topic which is beyond the scope of this study. Using more complex optimization algorithms instead of stochastic method to select target molecule would be an option for further study.

## 3 Results and Discussion

To evaluate t-SMILES and generatation models on different datasets, we firstly compare two t-SMILES based mGPT2 models with fragment-based baseline models and classical SMILES based mGPT2 models on Zinc. And then we compare t-SMILES based models with SOTA SMILES and GNNs baseline models on QM9 and ChEMBL datasets. Metrics used in this study are still based on distributional learning.

### 3.1 Comparison with fragment-based and SMILES-based models on Zinc

The validity, uniqueness, novelty, KLD and FCD scores on Zinc are summarized in Table 2 and the distributions of eight physiochemical properties are shown in Figure 3.

**Table 2.** Distributional results on Zinc, we train or retrained all these five models, t-SMILES_J_mGPT2 breaks molecules using the same fragmentation algorithm as JTVAE[27], t-Smiles_B_mGPT2 breaks molecules using BRICS, mGP2 means five layers mini GPT2 model is used.

| Model | Valid(↑) | Uniq(↑) | Novel(↑) | KLD(↑) | FCD(↑) |
|---|---|---|---|---|---|
| **SMILES_mGPT2(Ours)** | 0.853 | 0.674 | 0.672 | 0.960 | 0.830 |
| JTVAE[12] | 0.997 | 0.989 | 0.989 | 0.870 | 0.439 |
| FragDgm[27] | 1.000 | 0.423 | 0.422 | 0.835 | 0.303 |
| **t-SMILES_J_mGPT2(Ours)** | >0.99 | 0.775 | 0.774 | 0.970 | 0.773 |
| **t-Smiles_B_mGPT2(Ours)** | >0.99 | 0.783 | 0.782 | 0.963 | 0.790 |

Classical SMILES based mGPT2 model is trained for reference in this study to evaluate t-SMILES based models. Compared with t-SMILES models, classical SMILES based model gets lower novelty and uniqueness scores. It means that t-SMILES-based models, with almost 100% validity and relatively high FCD scores, could improve novelty to explore a wider molecular space.

JTVAE[12] is one key baseline model for this study, which splits molecule using a tree base data struct. So far as the validity is concerned, both t-SMILES and JTVAE models could generate near 100% valid molecules. However, t-SMILES exhibits significantly higher KLD and FCD scores with reasonably slight lower novelty and uniqueness scores.

t-SMILES based models significantly outperforms another fragment dictionary-based model, FragDgm[27], which splits molecule in a linear mode as a sequence of fragment IDs, on all five distribution parameters. Despite FragDgm adopting a segmented mode and based on distributional learning, its FCD value is the lowest among the five models.

When more and more validity of models can reach above 0.9, it becomes more important to test whether the generative model can effectively learn the physicochemical properties of molecules.

Detailed distributions of eight physiochemical properties logP, plogP, SAS, BertzCT, QED, TPSA, NPS and FractionCSP3 on Zin are presented in Figure 3.
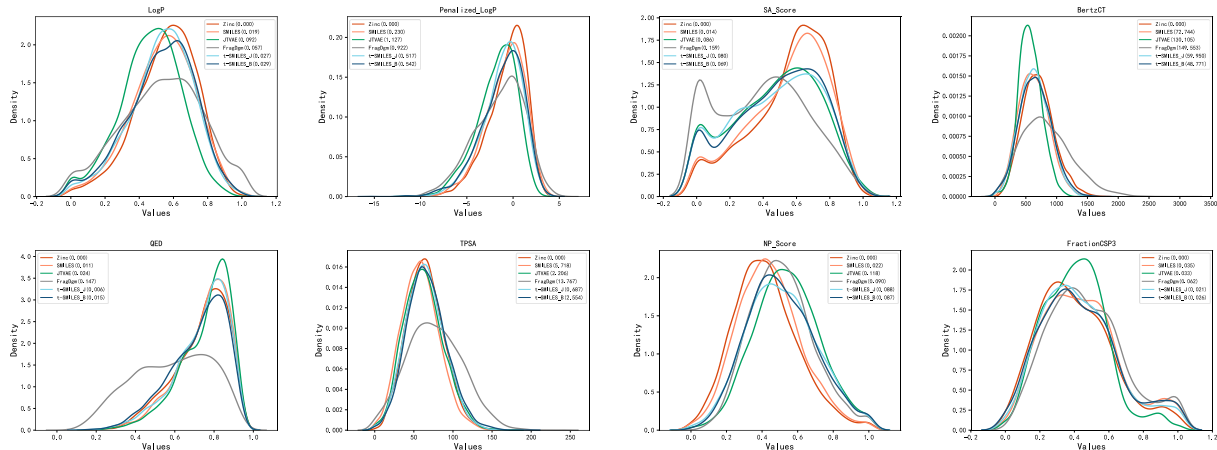


**Figure 3**. Distributions of physiochemical properties (logP, plogP, SAS, BertzCT, QED, TPSA, NPS and FractionCSP3) of random selected 10K molecules from Zinc data sets, JTVAE[12], FragDgm[27] and molecules generated by SMILES and two t-SMILES based mGPT2 models respectively

It can be seen from Figures 3 that the distribution of t-SMILES and classical SMILES based mGPT2 models are closer to training data set on all eight physiochemical properties than other two fragment-based models. The distribution of generated molecules by FragDgm is far from the training set on six physiochemical properties except NPS and FractionCSP3.

Compared with JTVAE, besides TPSA and SAS which are comparable, two t-SMILES based mGPT2 models get significantly lower scores on all eight metrics, especially on FractionCSP3 and BertzCT. Compared with FragDgm, except NPS which are comparable, t-SMILES based models get significantly favorite scores in other seven physicochemical properties.

Considering novelty, the distribution of the eight physicochemical properties suggests that the generated molecules by t-SMILES-based models are better fit to the training dataset, which obviously implies relatively low novelty.

From the comparative analysis of these models, two t-SMILES based models in this study significantly outperform the fragment-based baseline models, and have similar performance to the SMILES-based model in physicochemical properties, but validity scores could be greater than 99%, that is to say, t-SMILES could be an effective molecular presentation for fragment based molecular tasks on Zinc.

### 3.2 Comparison with baseline models on QM9

The validity, uniqueness, novelty, KLD and FCD scores on QM9 are summarized in Table 3 and the distributions of eight physiochemical properties are shown in Figure 4.

**Table 3.** Distributional results on QM9. The results of CharacterVAE[7], are taken from MolGAN[60], Transformer Reg[11], GraphVAE[59,11], MolGAN[60] and MGM[11] are taken from O. Mahmood *et al*.[11]

| Input | Model | Valid(↑) | Uniq(↑) | Novel(↑) | KLD(↑) | FCD(↑) |
|---|---|---|---|---|---|---|
| | CharacterVAE[7,60] | 0.103 | 0.675 | 0.900 | N/A | N/A |
| SMILES | Transformer Reg[11] | 0.965 | 0.957 | 0.183 | 0.994 | 0.958 |
| | **SMILES_mGPT2(ours)** | 0.949 | 0.728 | 0.172 | 0.992 | 0.970 |
| | GraphVAE[59,11] | 0.557 | 0.760 | 0.616 | N/A | N/A |
| Graph | MolGAN[60] | 0.981 | 0.104 | 0.942 | N/A | N/A |
| | MGM[11] | 0.886 | 0.978 | 0.518 | 0.966 | 0.842 |
| t-SMILES | **t-SMILES_J_mGPT2(ours)** | >0.99 | 0.720 | 0.289 | 0.976 | 0.953 |

On QM9, the score of validity of t-SMILES based mGPT2 model could be greater than 0.99, which is comparable to the performance of GNNs and superior to the most SOTA SMILES based models. FCD score of t-SMILES based mGPT2 model is one of the highest ones in all seven models, which is greater than 0.95.

From the viewpoint of sequence-based models, our approach performs similarly to or better than existing SMILES based approaches. Our approach shows higher validity and uniqueness scores compared to CharacterVAE, while having a reasonable lower novelty score. Compared to the Transformer Reg model, t-SMILES based model has higher score in novelty, lower score in uniqueness, and comparable scores in KLD and FCD.

Compared to the graph-based models, our approach outperforms the existing baseline approaches. t-SMILES based model has comparable uniqueness score compared with respect to GraphVAE, and significantly outperforms MolGAN, with reasonable lower novelty score. KLD and FCD scores are not provided for these two models. t-SMILES based mGPT2 model has good performance against SOTA method MGM in validity, KLD and FCD scores, while having slightly lower scores in uniqueness and novelty.
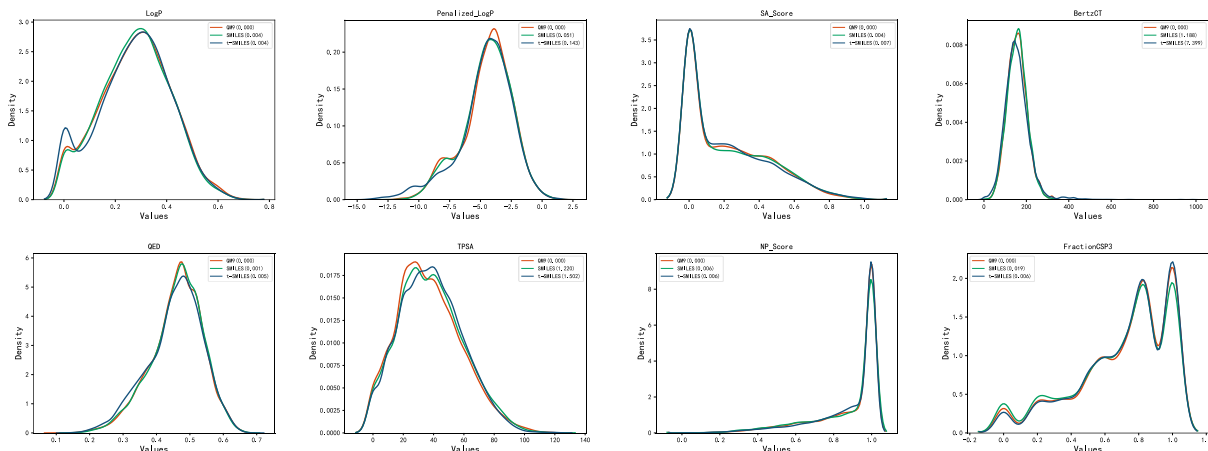


**Figure** 4. Distributions of physiochemical properties (logP, plogP, SAS, BertzCT, QED, TPSA, NPS and FractionCSP3) of random selected 10K molecules from generated molecules on QM9 data sets by SMILES and t-SMILES based mGPT2 models respectively

In general, on QM9 data set, both mGPT2 models based on t-SMILES and classical SMILES get high KLD and FCD scores, all greater than 0.95, that means they could maintain physiochemical similarity to the training distribution. However, the novelty scores of the models are relatively low, which could be interpreted by that the generated molecules based on distributional learning exhibit "the striking similarity" with known molecules included in the training set. It also indicates that novelty is inversely correlated with the KLD and FCD scores.

### 3.3 Comparison with baseline models on ChEMBL

The validity, uniqueness, novelty, KLD and FCD scores on ChEMBL are summarized in Table 4 and the distributions of eight physiochemical properties are shown in Figure 5.

**Table** 4 Distributional results on ChEMBL, The results of ORGAN[8], LSTM[37], CharacterVAE[7], AAE[61], Transformer Reg[11], Graph MCTS[14], and MGM[11] are taken from O. Mahmood *et al.*[11] and Guacamol[37], the results of MolGPT is taken from Bagal *et al.*[9], the results of hgraph2graph[13] is calculated by us.

| Input | Model | Valid(↑) | Uniq(↑) | Novel(↑) | KLD(↑) | FCD(↑) |
|-------|-------|----------|---------|----------|--------|--------|
| | ORGAN[8,37] | 0.379 | 0.841 | 0.687 | 0.267 | 0.000 |
| | LSTM[37] | 0.959 | 1 | 0.912 | 0.991 | 0.913 |
| | Character VAE[7,37] | 0.870 | 0.999 | 0.974 | 0.982 | 0.863 |
| SMILES | AAE[61,37] | 0.822 | 1 | 0.998 | 0.886 | 0.529 |
| | MolGPT[9] | 0.981 | 0.998 | 1 | 0.992 | 0.907 |
| | Transformer Reg[11] | 0.961 | 1.000 | 0.846 | 0.977 | 0.883 |
| | **SMILES_mGPT2(Ours)** | 0.850 | 0.670 | 0.641 | 0.972 | 0.809 |
| | Graph MCTS[14, 37] | 1.000 | 1.000 | 0.994 | 0.522 | 0.015 |
| Graph | MGM[11] | 0.849 | 1.000 | 0.722 | 0.987 | 0.845 |
| | hgraph2graph[13] | 1.000 | 0.994 | 0.940 | 0.870 | 0.485 |
| t-SMILES | **t-SMILES_J_mGPT2(Ours)** | >0.99 | 0.781 | 0.765 | 0.913 | 0.564 |
| | **t-SMILES_B_mGPT2(Ours)** | >0.99 | 0.781 | 0.769 | 0.935 | 0.575 |

On ChEMBL, two t-SMILES based mGPT2 models outperform graph-based baseline methods Graph MCTS and hgraph2graph, while the performance is comparable to SOTA graph mode MGM. Compared to Graph MCTS, t-SMILES based mGPT2 models show lower novelty scores while having significantly higher KLD and FCD scores. It seems difficult for this graph-based baseline model to capture the properties of the dataset distributions as shown by their low KLD scores and almost-zero FCD scores, but it gets the highest novelty score.

Compared to hgraph2graph which is an advanced model based on JTVAE and aims to solve larger molecular problems with motif-based method, t-SMILES based mGPT2 modes have higher KLD and FCD scores and lower uniqueness and novelty scores.

Compared to SOTA graph-based model MGM, t-SMILES based mGPT2 models have higher score in validity, but lower scores in uniqueness, KLD and FCD and similar score in novelty.

In the realm of sequence-based models, the proposed t-SMILES based mGPT2 models are competitive with the classical SMILES based models, and firstly, our models outperform all listed classical SMILES based models in validity. Compared with the GAN-based model (ORGAN), t-SMILES based mGPT2 models have significantly higher scores in validity, KLD and FCD and novelty scores, while having a slightly lower score in uniqueness. Compared with AAE, t-SMILES

modes get low unique and novelty scores and higher KLD and FCD scores. Similar to Graph MCTS, ORGAN also get an almost-zero FCD score. Our t-SMILES based approach results in lower scores across most of the metrics when compared to Transformer Reg, MolGPT, LSTM, VAE models besides validity score. In addition, statistics data shows that SMILES-mGPT2 model gets a slightly lower scores compared with Transformer Reg and MolGPT, which indicates that the t-SMILES mGPT2 models could be optimized to get high scores as well. Such an issue could be severed as a starting point for further research.
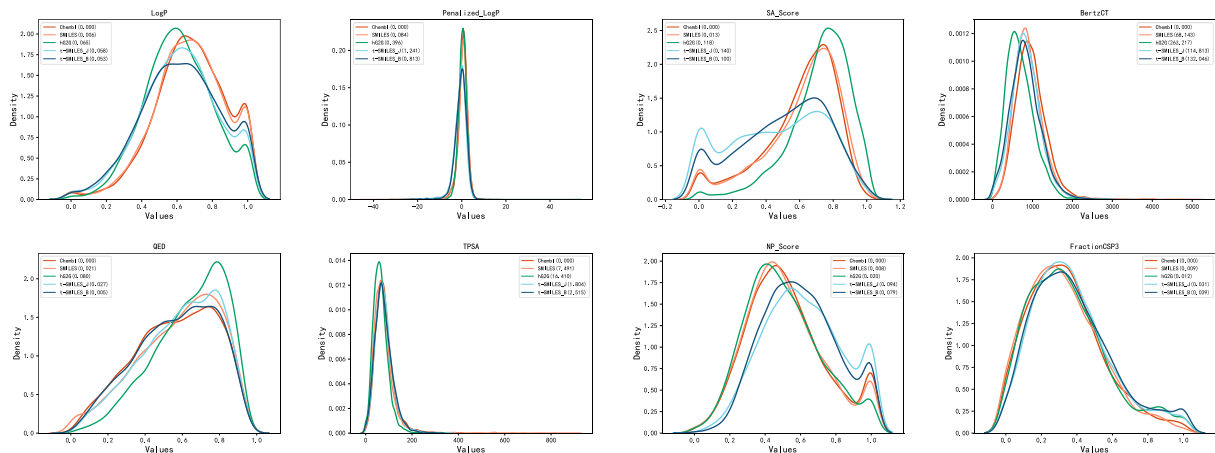


**Figure** 5. Distributions of physiochemical properties (logP, plogP, SAS, BertzCT, QED, TPSA, NPS and FractionCSP3) of random selected 10K molecules from generated molecules on ChEMBL data sets, hgraph2graph(hG2G)[13] and molecules generated by SMILES and t-SMILES based mGPT2 models respectively

Compared with hgraph2graph in Figure 5, t-SMILES based models get slightly higher scores on BertzCT, logP, QED and TPSA, while having comparable scores on SAS, NPS, plogP and FractionCSP3. Compared with classical SMILES based mGPT2 model, t-SMILES based mGPT2 models have lower scores on seven metrics besides QED.

On ChEMBL, one possible reason of the proposed t-SMILES models getting a relatively little lower scores might be that t-SMILES based mGPT2 model is not hyper-parametrically optimized well.

## 3.4 Ablation Studies

Compared to classical SMILES, only two more characters '&' and '^' are introduced for t-SMILES to encode multiscale molecular topological structure. Therefore, we adopt the same network structure to conduct ablation studies on one LSTM model and two Transformer-based models using three downstream datasets to investigate the performance of generative models. As shown in Table 5, different generative models get slightly different distributional results for the same dataset. However, compared to SMILES based models, t-SMILES based models get higher validity, uniqueness and novelty scores generally. This evidence also supports one of our claims that t-SMILES make it possible to generate molecules fragment-by-fragment using sequence-based models easily.

**Table** 5 Distributional results on Zinc, QM9 and ChEMBL datasets. The results of LSTM[37], MolGPT[9] and mGPT2 are all trained and calculated by us.

| Dataset | Model | Valid($\uparrow$) | Uniq($\uparrow$) | Novel($\uparrow$) | KLD($\uparrow$) | FCD($\uparrow$) |
|---|---|---|---|---|---|---|
| | t-SMILES_J_LSTM | >0.99 | 0.782 | 0.781 | 0.944 | 0.670 |
| Zinc | t-SMILES_J_MolGPT | >0.99 | 0.783 | 0.780 | 0.984 | 0.823 |
| | t-SMILES_J_mGPT2 | >0.99 | 0.775 | 0.774 | 0.970 | 0.773 |
| | t-SMILES_J_LSTM | >0.99 | 0.736 | 0.325 | 0.966 | 0.950 |
| QM9 | t-SMILES_J_MolGPT | >0.99 | 0.746 | 0.305 | 0.972 | 0.964 |
| | t-SMILES_J_mGPT2 | >0.99 | 0.720 | 0.289 | 0.976 | 0.953 |
| | t-SMILES_J_LSTM | >0.99 | 0.782 | 0.781 | 0.944 | 0.670 |
| ChEMBL | t-SMILES_J_MolGPT | >0.99 | 0.780 | 0.762 | 0.964 | 0.695 |
| | t-SMILES_J_mGPT2 | >0.99 | 0.781 | 0.765 | 0.913 | 0.564 |
| | SMILES_LSTM | 0.850 | 0.667 | 0.666 | 0.966 | 0.831 |
| Zinc | SMILES_MolGPT | 0.980 | 0.770 | 0.766 | 0.989 | 0.923 |
| | SMILES_mGPT2 | 0.853 | 0.674 | 0.672 | 0.960 | 0.830 |
| | SMILES_LSTM | 0.899 | 0.690 | 0.182 | 0.992 | 0.965 |
| QM9 | SMILES_MolGPT | 0.965 | 0.729 | 0.150 | 0.986 | 0.907 |
| | SMILES_mGPT2 | 0.949 | 0.728 | 0.172 | 0.992 | 0.970 |
| | SMILES_LSTM | 0.774 | 0.609 | 0.598 | 0.943 | 0.725 |
| ChEMBL | SMILES_MolGPT | 0.976 | 0.768 | 0.711 | 0.992 | 0.807 |
| | SMILES_mGPT2 | 0.850 | 0.670 | 0.641 | 0.972 | 0.809 |

## 4 Experiments

In this section, we briefly introduce the proposed generation models which are used to evaluate t-SMILES. This is then followed by the overview of three commonly used public data sets. Finally,

the details of the experiments and the metrics used for the evaluation of different models are provided.

## 4.1 Generation models and hyperparametric optimization

Language modelling uses probability and statistical techniques to calculate the possibility of word sequences in sentences and then do estimation[62]. Originally, recurrent neural network (RNN) is designed to address this kind of sequence problems. To address the limitation of RNN that lead to gradient disappearance and explosion for long sequences, cyclic models such as LSTM and GRU are proposed. However, it has been shown that the power of LSTMs is insufficient when the information has ultra-long dependencies. In 2017, the attention-based Transformer architecture broke through the limitations of LSTM and quickly achieved SOTA scores on multiple metrics in NLP and computer vision. After that, the improved models such as GPT1-3[63,43,64]and BERT[65] are introduced to build pre-trained model on large datasets. Although these large-size models have achieved unprecedented performances, they come at high computational costs. Consequently, some of the recent NLP architectures have utilized concepts of transfer learning, pruning, quantization, and knowledge distillation to achieve moderate model sizes while keeping nearly similar performances as achieved by their predecessors. The recent developments in the NLP field have a great potential to be adapted to molecular de novo generation research. Bagal *et al.*[9] has proposed a Transformer-decoder model with 10 layers to molecular generation task and proved it works well.

In this study, we mainly adopt Transformer-decoder based autoregressive generation models to evaluate our proposed t-SMILES. For the sake of data comparability, we select one model as the basis using a Bayesian optimizer SigOpt[66] to optimize the hyperparameters related to the neural network topology in ChEMBL dataset using classical SMILES as input. This selected model is a

5 layers mini version of GPT2(mGPT2) with 512 hidden layer and 8 attention headers. All other mGPT2 models in different dataset with either t-SMILES or classical SMILES use the same high-level neural network architecture. Finally, we train all models in detail using the Adam optimizer.

4.2 Datasets

We evaluate t-SMILES on three commonly used public data sets including subset of Zinc[32], QM9[33] and ChEMBL-21[34].

The Zinc subset used in our experiments is the same as the subset in JTVAE which contains approximately 250K drug-like molecules extracted from the Zinc15 database that span a wide range of chemistry. The molecules in this Zinc subset are composed of 10 atoms including Br，C，Cl，F，H，I，N，O，P and S.

QM9 dataset contains up to 9 heavy (non hydrogen) atoms. Molecules in the subset consist of 5 atoms including C, F, N, H and O atoms. In all, this results in about 134k druglike organic molecules.

For the sake of universality, we select a subset from ChEMBL-21 with SMILES character length less than 120. Apart from this very simple rule, no other preprocessing is done. Molecules in ChEMBL subset are composed of 11 atoms including B, Br, C, Cl, F, H, I, N, O, P and S.

4.3 Evaluation Metrics

Despite a large number of metrics have been proposed for the evaluation of generative models of molecules, there is little consensus on which should be used. It is often biased to use simple indexes to evaluate the performance of different models [67]. From the perspective of optimization, it could be considered that the task can be solved once the molecules with a high score of a certain index are generated, but these generated molecules may not be useful. Therefore, we use five benchmarks

proposed by GuacaMol: validity (↑), uniqueness (↑), novelty (↑), KL divergence (KLD) (↑)[35]and

Fréchet ChemNet Distance(FCD)[36] score (↑) to evaluate the general performance of the model.

GuacaMol calculate every score using 10K randomly sampled molecules. Validity measures the

ratio of valid molecules which could be correctly parsed by RDKit[45]; Uniqueness is the

percentage of valid molecules that are unique; Novelty is the percentage of valid unique molecules

that are not included in the training data set; KLD[35] score compares the distribution of a variety

physicochemical descriptors of the training set and generated molecules. FCD[36] score measures

the proximity of the distribution of generated molecules to the distribution of the dataset molecules

according to the Fréchet Distance in the hidden representation space of ChemNet[68], which is

trained to predict the chemical properties of small molecules. The values of these five parameters

are between 0 and 1. The larger the value, the 'better' the model. While KLD and FCD are both

measure of the similarity between generated molecules and molecules from the training data set.

They are highly correlated with each other and inversely correlated with novelty[11].

In addition, we use logP, plogP[38], SAS[39], BertzCT[40], QED[41], TPSA, NPS[42], and

FractionCSP3 to evaluate whether the generative models could effectively learn the physical and

chemical properties of the molecules in the training set, thereby comprehensively evaluating the

performance of the generative model from the perspective of distributed learning knowledge.

- **logP**: The logarithm of the partition coefficient. If one of the solvents is water and the other is a nonpolar solvent, then logP is a measure of hydrophobicity.

- **Penalized logP(plogP**)[38]is the logarithm of the partition ratio of solute between octanol and water subtracted by synthetic accessibility score and long cycles. It has a range of $(-\infty, \infty)$

- **Synthetic Accessibility score (SAS)**[39]: Measurement of the difficulty of synthesizing a compound. It is a score between 0 (difficult) and 1(easy).

- **BertzCT**[40]: A topological index meant to quantify "complexity" of molecules. It consists of a sum of two terms, one representing the complexity of the bonding, the other representing the complexity of the distribution of heteroatoms.

- **Quantitative Estimate of Drug-likeness (QED)**[41]: This quantifies drug-likeness by considering the main molecular properties. It ranges from 0 (all properties unfavorable) to 1 (all properties favorable).

- **Topological Polar Surface Area (TPSA)**: The sum of surface area over all polar atoms. It measures the drug's ability to permeate cell membranes. Molecules with a TPSA greater than 140 $\text{Å}^2$ tend to be poor in permeating cell membranes.

- **Natural Product-likeness score (NPS)**[42]: Score of similarity degree of structural space covered by molecules and natural products. It has a range of (0,1)

- **FractionCSP3**: The fraction of C atoms that are SP3 hybridized.

In general, we randomly select 10K molecules from generated ones to calculate five metrics from GuacaMol[37] and the distribution of eight physiochemical properties.

4.4 Details of baseline models

We train classical SMILES based mGPT2 models as relatively baseline on three datasets, which share the same neural network architecture with t-SMILES. If there are some SMILS-based models that get better scores than this SMILES-mGPT2, it means that it's possible to build t-SMILES based model to get better scores.

In ablation studies, we retrain the LSTM[37] model with 3 layers of hidden size 768, dropout of 0.2 using the Adam optimizer with learning rate 0.001, and MolGPT[9] model with 8 layers, 8

headers with hidden size 256 using AdamW optimizer with learning rate 0.001 for both SMILES and t-SMILES on all three datasets.

We retrain JTVAE[12] on Zinc using the publicly available codebase provided by the paper's authors and then generate 20K molecules.

We calculate metrics based on the molecules provided by Fragment-base-DGM[27] on Zinc.

We generate 20K molecules using the pre-trained model provided by hgraph2graph[13] on ChEMBL and calculate metrics for evaluation.

We do not train the rest of the baseline models by ourselves. For QM9 and ChEMBL, we take some results of baseline models from Bagal *et al.*[9], Cao and Kipf [60], O. Mahmood *et al.*[11] and GuacaMol[37].

Finally, we follow the open-source implementation of the GuacaMol[37] benchmark baselines to calculate metrics for all baseline and new designed models.

## 5 Conclusion and Outlook

One of the challenging problems in designing drug molecules is the vast chemical space to be explored. Fragment-based approaches identify a subset of chemical moieties responsible for key molecular recognitions early on and this allows scientists to devote their time in a much reduced and relevant chemical space.

In general, if atoms are analogized to characters in natural language processing, fragments could be analogized to words or phrases as functional "units", then molecule could be a fragment based discrete structured data. To train a general-purpose string generative mode to learn the knowledge from discrete structured data, we need to prepare large amount of valid combinations of the structures, which is time consuming and will also face the challenge of insufficient effective data

practical in domains like target-oriented drug discovery. Our proposed t-SMILES uses SMILES instead of fragment dictionaries id, re-encodes fragment-based molecular tree through a full binary tree and generates molecular sequence with multiscale structural information, thus providing a new idea for fragment-based molecular design. In this way, powerful and rapidly developing sequence-based solutions can be applied to fragment-based molecular tasks in the same way as classical SMILES.

In addition, compared with classical SMILES which is relatively difficult to be augmented, by using different fragmentation algorithms, the training dataset is easier and more efficiently to be expanded on t-SMILES to explore different chemical spaces without having to change anything of the architecture of generation model.

**Scalability** Our proposed t-SMILES solution supports any effective substructures types and patterns as long as they could be used to obtain chemically valid molecular fragments and ultimately construct valid acyclic molecular trees. So that, with the invaluable accumulation of drug fragments by experienced chemists, it's possible that the t-SMILES algorithm combines this experience with a powerful sequence based deep neural network model to help chemists better explore chemical space. If in some special cases, only one specific fragment space needs to be explored, and expansion is not required, then in t-SMILES algorithm, tree nodes could be easily encoded by dictionary id instead of SMILES. Or we could replace the newly generated fragments with the fragments from the training data set according to specified rules. The encoding logic and algorithm flow of t-SMILES remain unchanged.

**To be improved** Our experiments show that how to segment, assemble molecular fragments and how to optimize molecular are also key steps controlling the quality of generated molecules. Therefore, this challenging topic could serve as a starting point for further research.

With the rapid development of NLP technology, Transform-based models have been proved to enable text generation with human-like capabilities if trained on enough data. SMILES-based models have proved that the sequence-based NLP models are powerful tools to generate molecules. Then, t-SMILES makes it easily to use sequence-based NLP models to generate molecules fragment-by-fragment. Therefore, the investigation of t-SMILES based generation models or other fragment-based tasks would be served as interesting topic for further study.

## List of abbreviations

SOTA: State-of-the-art

FBDD: Fragment-based drug discovery

AMT: Acyclic Molecular Tree

FBT: Full Binary Tree

BFS: Breadth First Search

DFS: Depth First Search

SMILES: Simplified Molecular Input Line Entry Specification

plogP: Penalized logP

SAS: Synthetic Accessibility score

QED: Quantitative Estimate of Drug-likeness

NPS: Natural Product-likeness score

FCD: ChemNet distance score

KLD: KL divergence

## Conflicts of Interest

The authors declare that they have no competing interests.

## Author Contributions

Ruqin Yu and Juanni Wu designed the study and manuscript. Juanni Wu conceived the project, constructed the algorithms and Python script, performed the experiments, informatics analyses, and wrote the draft manuscript. Tong Wang and Yue Chen participated in the discussion and experiments. Lijuan Tang and Hailong Wu participated in the discussion and funding acquisition. All authors contributed to manuscript editing, revising and have approved the final version of the manuscript.

## Data availability

The datasets used in this study are publicly available. They are referenced in the Datasets and Evaluation Metrics part of the Experiments section. The processed data used in this study can be found at: https://github.com/juanniwu/t-SMILES/

## Code availability

Code, pretrained t-SMILES models, training and generation scripts for this work and lists of generated molecules can be found at: https://github.com/juanniwu/t-SMILES/

The code of baseline models used to in this work are publicly available. We are gratefully acknowledging all authors of these researches.

[1] MolGPT[9]: https://github.com/devalab/molgpt

[2] MGM[11]: https://github.com/nyu-dl/dl4chem-mgm

[3] JTVAE[12]： https://github.com/wengon-jin/icml18-jtnn

[4] hgraph2graph[13]: https://github.com/wengong-jin/hgraph2graph

[5] FragDGM[27]: https://github.com/marcopodda/fragment-based-dgm

[6] Guacamol[37]：https://github.com/BenevolentAI/guacamol_baselines

[7] GPT2: https://github.com/samwisegamjeee/pytorch-transformers

**Reference**

1.    Shuker SB, Hajduk P, Meadows R, Fesik S (1996) Discovering High-Affinity Ligands for
      Proteins: SAR by NMR. Science (80- ) 274:1531–1534.
      https://doi.org/10.1126/science.274.5292.1531

2.    Erlanson DA, Fesik SW, Hubbard RE, et al (2016) Twenty years on: The impact of
      fragments on drug discovery. Nat Rev Drug Discov 15:605–619.
      https://doi.org/10.1038/nrd.2016.109

3.    Lecun Y, Bengio Y, Hinton G (2015) Deep learning. Nature. 521:436–444

4.    Butler K, Davies DW, Cartwright H, et al (2018) Machine learning for molecular and
      materials science. Nature, 2018, 559(7715): 547-555.

5.    Schneider P, Walters WP, Plowright A, et al (2019) Rethinking drug design in the
      artificial intelligence era. Nat Rev Drug Discov 19:353–364.
      https://doi.org/10.1038/s41573-019-0050-3

6.    Weininger D (1988) SMILES, a Chemical Language and Information System: 1:
      Introduction to Methodology and Encoding Rules. J Chem Inf Comput Sci 28:31–36.
      https://doi.org/10.1021/ci00057a005

7.    Gómez-Bombarelli R, Wei JN, Duvenaud D, et al (2018) Automatic Chemical Design

      Using a Data-Driven Continuous Representation of Molecules. ACS Cent Sci 4:268–276.

      https://doi.org/10.1021/acscentsci.7b00572

8.    Guimaraes GL, Sánchez-Lengeling B, Farias PLC, Aspuru-Guzik A (2017) Objective-

      Reinforced Generative Adversarial Networks (ORGAN) for Sequence Generation Models.

      ArXiv:1705.10843, 2017

9.    Bagal V, Aggarwal R, Vinod PK, Priyakumar UD (2021) MolGPT: Molecular Generation

      Using a Transformer-Decoder Model. J Chem Inf Model. 2021, 62(9): 2064-2076

10.   Xia X, Hu J, Wang Y, et al (2019) Graph-based generative models for de Novo drug

      design. Drug Discov Today Technol 32–33:45–53.2019

      https://doi.org/10.1016/j.ddtec.2020.11.004

11.   Mahmood O, Mansimov E, Bonneau R, Cho K (2021) Masked graph modeling for

      molecule generation. Nat Commun 12:1–36. https://doi.org/10.1038/s41467-021-23415-2

12.   Jin W, Barzilay R, Jaakkola T (2018) Junction Tree Variational Autoencoder for

      Molecular Graph Generation. In: International conference on machine learning. PMLR, pp

      2323–2332

13.   Hu Y, Hu Y, Cen E (2021) Hierarchical Generation of Molecular Graphs using Structural

      Motifs. Proc - 2021 2nd Int Conf Big Data Artif Intell Softw Eng ICBASE 2021 543–546.

      https://doi.org/10.1109/ICBASE53849.2021.00106

14.   Jensen JH (2019) A graph-based genetic algorithm and generative model/Monte Carlo tree

      search for the exploration of chemical space. Chem Sci 10:3567–3572.

      https://doi.org/10.1039/c8sc05372c

15.	Hoogeboom E, Satorras VG, Vignac C, Welling M (2022) Equivariant Diffusion for Molecule Generation in 3D. ArXiv.2022

16.	Bodnar C, Frasca F, Otter N, et al (2021) Weisfeiler and Lehman go cellular: CW networks. Adv Neural Inf Process Syst 34:2625–2640

17.	Bouritsas G, Frasca F, Zafeiriou S, Bronstein M (2020) Improving Graph Neural Network Expressivity via Subgraph Isomorphism Counting. ArXiv abs/2006.0: https://doi.org/10.1109/TPAMI.2022.3154319

18.	Bodnar C, Frasca F, Wang Y, et al (2021) Weisfeiler and lehman go topological: Message passing simplicial networks. In: International Conference on Machine Learning. PMLR, pp 1026–1037

19.	Wu Z, Pan S, Chen F, et al (2021) A Comprehensive Survey on Graph Neural Networks. IEEE Trans Neural Networks Learn Syst 32:4–24. https://doi.org/10.1109/TNNLS.2020.2978386

20.	Skinnider M, Stacey RG, Wishart D, Foster L (2021) Chemical language models enable navigation in sparsely populated chemical space. Nat Mach Intell 3:759–770. https://doi.org/10.1038/S42256-021-00368-1

21.	O'Boyle NM, Dalke A (2018) DeepSMILES: An adaptation of SMILES for use in machine-learning of chemical structures. ChemRxiv 1–9. https://doi.org/10.26434/CHEMRXIV.7097960.V1

22.	Krenn M, Häse F, Nigam A, et al (2019) Self-Referencing Embedded Strings (SELFIES): A 100% robust molecular string representation. Mach Learn Sci Technol 1:045024. https://doi.org/10.1088/2632-2153/aba947

23.    Mitrovic J, McWilliams B, Walker J, et al (2020) Representation Learning via Invariant Causal Mechanisms. ICML 1–21

24.    Maziarz K, Jackson-Flux H, Cameron P, et al (2021) Learning to Extend Molecular Scaffolds with Structural Motifs. ArXiv abs/2103.0

25.    Zhang Z, Liu Q, Wang H, et al (2021) Motif-based Graph Self-Supervised Learning for Molecular Property Prediction. ArXiv abs/2110.0

26.    Yu Z, Gao H (2022) Molecular Representation Learning via Heterogeneous Motif Graph Neural Networks. ICML

27.    Podda M, Bacciu D, Micheli A (2020) A deep generative model for fragment-based molecule generation. In: International Conference on Artificial Intelligence and Statistics. PMLR, pp 2240–2250

28.    Jin W, Barzilay R, Jaakkola T (2020) Multi-Objective Molecule Generation using Interpretable Substructures. 37th Int Conf Mach Learn ICML 2020 PartF16814:4799–4809

29.    Vaswani A, Shazeer N, Parmar N, et al (2017) Attention is all you need. Adv Neural Inf Process Syst 2017-Decem:5999–6009

30.    Flam-Shepherd D, Zhu K, Aspuru-Guzik A (2022) Language models can learn complex molecular distributions. Nat Commun 13:1–10. https://doi.org/10.1038/s41467-022-30839-x

31.    Degen J, Wegscheid-Gerlach C, Zaliani A, Rarey M (2008) On the art of compiling and using "drug-like" chemical fragment spaces. ChemMedChem 3:1503–1507. https://doi.org/10.1002/cmdc.200800178

32.  Sterling T, Irwin JJ (2015) ZINC 15 - Ligand Discovery for Everyone. J Chem Inf Model 55:2324–2337. https://doi.org/10.1021/acs.jcim.5b00559

33.  Ramakrishnan R, Dral PO, Rupp M, Von Lilienfeld OA (2014) Quantum chemistry structures and properties of 134 kilo molecules. Sci Data 1:1–7. https://doi.org/10.1038/sdata.2014.22

34.  Gaulton A, Bellis LJ, Bento AP, et al (2012) ChEMBL: A large-scale bioactivity database for drug discovery. Nucleic Acids Res 40:. https://doi.org/10.1093/NAR/GKR777

35.  Kullback S, Leibler RA (1951) On Information and Sufficiency. Ann Math Stat 22:79–86. https://doi.org/10.1214/AOMS/1177729694

36.  Preuer K, Renz P, Unterthiner T, et al (2018) Fréchet ChemNet Distance: A Metric for Generative Models for Molecules in Drug Discovery. J Chem Inf Model 58:1736–1741. https://doi.org/10.1021/acs.jcim.8b00234

37.  Brown N, Fiscato M, Segler MHS, Vaucher A (2019) GuacaMol: Benchmarking Models for De Novo Molecular Design. J Chem Inf Model 59 3:1096–1108. https://doi.org/10.1021/acs.jcim.8b00839

38.  Rajasekar AA, Raman K, Ravindran B (2020) Goal directed molecule generation using Monte Carlo Tree Search. ArXiv abs/2010.1:

39.  Ertl P, Schuffenhauer A (2009) Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions. J Cheminform 1:. https://doi.org/10.1186/1758-2946-1-8

40.  Bertz SH (1981) The First General Index of Molecular Complexity. J Am Chem Soc 103:3599–3601.

https://doi.org/10.1021/JA00402A071/ASSET/JA00402A071.FP.PNG_V03

41.  Bickerton GR, Paolini G V., Besnard J, et al (2012) Quantifying the chemical beauty of drugs. Nat Chem 4:90–98. https://doi.org/10.1038/NCHEM.1243

42.  Ertl P, Roggo S, Schuffenhauer A (2008) Natural product-likeness score and its application for prioritization of compound libraries. J Chem Inf Model 48:68–74. https://doi.org/10.1021/ci700286x

43.  Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei IS (2020) Language Models are Unsupervised Multitask Learners. OpenAI 1:1–7

44.  Lounkine E, Batista J, Bajorath J (2008) Random Molecular Fragment Methods in Computational Medicinal Chemistry. Curr Med Chem 15:2108–2121. https://doi.org/10.2174/092986708785747607

45.  Landrum G (2013) RDKit: A software suite for cheminformatics, computational chemistry, and predictive modeling. Greg Landrum

46.  Lewell XQ, Judd DB, Watson SP, Hann MM (1998) RECAP - Retrosynthetic Combinatorial Analysis Procedure: A powerful new technique for identifying privileged molecular fragments with useful applications in combinatorial chemistry. J Chem Inf Comput Sci 38:511–522. https://doi.org/10.1021/ci970429i

47.  Rarey M, Dixon JS (1998) Feature trees: A new molecular similarity measure based on tree matching. J Comput Aided Mol Des 12:471–490. https://doi.org/10.1023/A:1008068904628

48.  Clayden J, Greeves N, Warren SG (2001) Organic chemistry. Oxford University Press

49.  Takahashi Y, Sukekawa M, Sasaki S ichi (1992) Automatic Identification of Molecular

Similarity Using Reduced-Graph Representation of Chemical Structure. J Chem Inf Comput Sci 32:639–643. https://doi.org/10.1021/ci00010a009

50. Rarey M, Stahl M (2001) Similarity searching in large combinatorial chemistry spaces. J Comput Aided Mol Des 15:497–520. https://doi.org/10.1023/A:1011144622059

51. Leach AR, Gillet VJ (2007) An introduction to chemoinformatics. Springer

52. Gao W, Mercado R, Coley CW (2021) Amortized Tree Generation for Bottom-up Synthesis Planning and Synthesizable Molecular Design. arXiv Prepr arXiv211006389

53. Liu Y, Mathis C, Bajczyk MD, et al (2021) Exploring and mapping chemical space with molecular assembly trees. Sci Adv 7:. https://doi.org/10.1126/sciadv.abj2465

54. Bradshaw J, Paige B, Kusner MJ, et al (2020) Barking up the right tree: an approach to search over molecule synthesis DAGs. ArXiv abs/2012.1: https://doi.org/10.17863/CAM.69874

55. Nguyen DH, Tsuda K (2021) A generative model for molecule generation based on chemical reaction trees. 1–12

56. Ertl P, Schuffenhauer A, Renner S (2011) Reduced Graphs and Their Applications in Chemoinformatics. Life Sci 672:588. https://doi.org/10.1007/978-1-60761-839-3

57. Koller D, Friedman N (2009) Probabilistic graphical models: principles and techniques. MIT press

58. Mordatch I (2018) MOG: MOLECULAR OUT-OF-DISTRIBUTION GENERATION WITH ENERGY-BASED MODELS. 1–6

59. Simonovsky M, Komodakis N (2018) GraphVAE: Towards generation of small graphs using variational autoencoders. Lect Notes Comput Sci (including Subser Lect Notes Artif

Intell Lect Notes Bioinformatics) 11139 LNCS:412–422. https://doi.org/10.1007/978-3-030-01418-6_41

60. De Cao N, Kipf T (2018) MolGAN: An implicit generative model for small molecular graphs. ArXiv

61 Makhzani A, Shlens J, Jaitly N, et al (2015) Adversarial Autoencoders. arXiv, abs/1511.05644 (2015).

62. Singh S, Mahmood A (2021) The NLP Cookbook: Modern Recipes for Transformer Based Deep Learning Architectures. IEEE Access 9:68675–68702. https://doi.org/10.1109/ACCESS.2021.3077350

63. Radford A, Narasimhan K, Salimans T, Sutskever I (2018) Improving Language Understanding by Generative Pre-Training. Homol Homotopy Appl

64. Brown TB, Mann B, Ryder N, et al (2020) Language models are few-shot learners. Adv Neural Inf Process Syst 2020-Decem:

65. Devlin J, Chang MW, Lee K, Toutanova K (2019) BERT: Pre-training of deep bidirectional transformers for language understanding. NAACL HLT 2019 - 2019 Conf North Am Chapter Assoc Comput Linguist Hum Lang Technol - Proc Conf 1:4171–4186

66. Dewancker I, McCourt M, Clark S, et al (2016) Evaluation System for a Bayesian Optimization Service. arXiv Prepr

67. Renz P, Van Rompaey D, Wegner JK, et al (2019) On failure modes in molecule generation and optimization. Drug Discov Today Technol 32–33:55–63. https://doi.org/10.1016/j.ddtec.2020.09.003

68. Goh GB, Siegel C, Vishnu A, Hodas NO (2017) ChemNet: A Transferable and

Generalizable Deep Neural Network for Small-Molecule Property Prediction. arXiv

abs/1712.0:1712.02734