# Comparative Study of Deep Generative Models on Chemical Space Coverage

*Jie Zhang[&,$,||], Rocío Mercado[∈], Ola Engkvist[∈], Hongming Chen[&,||],\**

[&]Guangdong Provincial Key Laboratory of Laboratory Animals, Guangdong Laboratory

Animals Monitoring Institute, Guangzhou, 510663, P. R. China

[$]State Key Laboratory of Respiratory Disease, Guangzhou Institutes of Biomedicine and Health,

Chinese Academy of Sciences, Guangzhou 510530, P. R. China

[||]Bioland Laboratory (Guangzhou Regenerative Medicine and Health - Guangdong

Laboratory)，Guangzhou 510530, P. R. China

[∈]Hit Discovery, Discovery Sciences, R&D, AstraZeneca, Gothenburg 43183, Sweden

[\*]Correspondence e-mail: chen_hongming@grmh-gdl.cn

## Abstract

In recent years, deep molecular generative models have emerged as novel methods for *de novo* molecular design. Thanks to the rapid advance of deep learning techniques, deep learning architectures such as recurrent neural networks, variational autoencoders, and adversarial networks, have been employed for constructing generative models. However, so far the metrics used to evaluate these deep generative models are not discriminative enough to separate the performance of various state-of-the-art generative models. This work presents a novel metric for evaluating deep molecular generative models; this new metric is based on the chemical space coverage of a

reference database, and compares not only the molecular structures, but also the ring systems and

functional groups, reproduced from a reference dataset of a 1M subset of GDB-13. The

performance of 7 different molecular generative models was compared by calculating their

structure and substructure coverage of the GDB-13 database while using the 1M subset for training.

The result shows that the performance of various generative models varies significantly using the

benchmarking metrics introduced herein, such that generalization capability of the generative

model can be clearly differentiated. Additionally, the coverage of ring systems and functional

groups existing in GDB-13 was also compared between the models. Our study provides a useful

new metric that can be used for evaluating and comparing generative models.

## **Introduction**

Deep learning has been successfully used in many fields to learn relationships that are too complex

to learn using traditional computer algorithms, including early image classification,[1, 2] facial

recognition, and music recognition.[3] Deep learning even surpasses the performance of human

experts in some challenging tasks, such as playing GO.[4] Moreover, deep generative models play

important roles in tasks like music composition,[5] image generation,[6] and language translation. [7] In

the last five years, deep generative modeling has also been applied in the fields of cheminformatics

and molecular design. One interesting example is using deep neural networks for compound

structure generation.[8-11]

The number of chemically feasible, drug-like molecules has been estimated to be on the order of

$10^{60} - 10^{100}$ compounds.[12] For such a large chemical space, it is clearly impossible to synthesize and

test every compound for pharmaceutical applications. To efficiently explore the space, molecular

generative models have emerged in recent years with the aim of better navigating through this huge chemical space for *de novo* molecule design.

*De novo* molecular design has long been put forward as a way to accelerate the drug discovery process as it is expected to save time and resources in drug discovery, where it can take over a decade and billions of dollars in investment to bring a single drug to market.[13] Historically, *de novo* design methods have been mainly rule-driven and used brute force algorithms to achieve their goal.[14] For example, creating a virtual library using fixed rules and building blocks, then scoring each compound in the virtual library to find the best compound. Genetic algorithm based algorithms were also proposed to tackle the *de novo* design issue.[15,16] In contrast, deep generative molecular design is the concept of generating molecules using deep neural networks. Deep generative models are data-driven methods which generate compound structures by learning the underlying probability distributions in a compound dataset instead of screening existing databases for molecules that fit the desired profile. Deep generative models are powerful as they allow chemists to bypass models using hard-coded chemical rules which do not scale to larger datasets. Furthermore, not all chemical rules are easy to define. Using deep generative models, one can avoid enumerating all possible structures for a given application and then screening them (a daunting task). Instead, one can simply train a model using known compounds, and sample the model for the desired set of properties (e.g. ADMET profile) to get out promising structures. *De novo* generative models can generate structures that are in significantly narrower, but more promising, regions of chemical space. Moreover, deep learning methods can take advantage of all the information available in ever-increasing large public datasets, thanks to automation technologies used in high-throughput screening and parallel synthesis.[17]

In recent years, many molecular generative models have been published, such as CharRNN, VAE, and REINVENT, which are remarkable at sampling molecules both in- and outside the training set used to learn chemistry rules.[11,18-21] It is worth noting that CharRNN was introduced as a general language model at the first place. However, similar architectures are also successfully applied in molecular generative models, e.g. REINVENT adopted a similar architecture with reinforcement learning.[22-26] VAE is a general architecture that has a wide range of applications in many generative models and tasks.[27-29] In current study, we adopts the implementation of CharRNN and VAE provided by the MOSES.[10] Notably, many of these generative models have been benchmarked using existing "distribution-based" metrics implemented in open-source programs such as MOSES[10] or GuacaMol.[30] However, these metrics are in general non-discriminative as many of these state-of-the-art (SOTA) models perform quite well across all the included metrics, such that it is difficult to compare them and gain a deep understanding of each model's strengths and weaknesses. We previously proposed a new metric: the percent coverage of functional groups present in GDB-13.[31]  As an extension of our previous work,[32] we apply the idea as a way for benchmarking the performance of multiple generative models. GDB-13 contains in total 975,820,210 structures, which enumerate small organic molecules containing up to 13 atoms of C, N, O, S, and Cl by following simple chemical stability and synthetic feasibility rules.[32] The generalization capability of deep generative models is assessed by computing how much of the whole GDB-13 can be recovered by a model trained from a small GDB-13 subset.

Substructure has long been used to characterize the composition of compounds. One concept is the so-called *functional group*, frequently used in many fields in chemistry, including medicinal chemistry. A functional group is defined as a subset of connected atoms in a molecule that in some way endows specific intrinsic properties (or *functions*) to a molecule. Furthermore, the presence

86  or absence of a functional group in a molecule could determine whether a molecule will react in a

87  given reaction. Some of the most common groups in medicinal chemistry include amides

88  (RC(=O)NR'R''), ethers (R–O–R'), and amines (RR'NR''), where R, R', and R'' represent organic

89  groups or hydrogen atoms.[33] Another substructure-based concept is the *ring system*; ring systems

90  are the key components of molecular scaffolds. They play an important role in a molecule's

91  observed properties, such as the electronics, scaffold rigidity, molecular reactivity, and toxicity.

92  On average six new ring systems enter the drug space each year and approximately 28% of new

93  drugs contain a new ring system.[34] We investigated the percentage of chemical space covered in

94  terms of full structures, functional groups, and ring systems by published SOTA generative models.

95  The size of GDB-13 was hypothesized to be large enough to highlight differences between the

96  various models.

97  Four major classes of deep generative models are benchmarked and studied in this work, including

98  those based on recurrent neural networks (RNNs), autoencoder (AE) based networks, generative

99  adversarial networks (GANs), and graph neural networks (GNNs). The deep generative models

100 based on RNNs include REINVENT[18, 32, 35] and CharRNN,[25] which use SMILES as the input and

101 output strings. VAE,[36] AAE,[21,37] ORGAN,[20] and LatentGAN[11] adopt either an AE or GAN for structure

102 generation using SMILES. Besides the SMILES-based generative models, one graph-based

103 generative model, GraphINVENT,[38] which uses GNNs in its core architecture, is also included in

104 the benchmark study. An effort was made to cover most of the major types of generative model

105 architectures in this study, in the hope that this would provide a comprehensive comparison for

106 existing generative models.

## Methods

**Extraction of functional groups and ring systems.** To identify functional groups (FG) in the various sets of molecules in this work (generated molecule sets, and GDB-13), the RDKit functional group identification package,[39] which is based on an algorithm introduced by Peter Ertl for automatically identifying functional groups, was used.[40] The advantage of the method is that it is not based on manually curated lists of functional groups, and thus can be applied to any chemical series. It is important to note that different chemists have slightly different definitions of what is a functional group; however, as the benchmark introduced here calculated ratios of functional groups in the generated and reference sets, a difference in opinions between chemists shouldn't be relevant. The extraction of compound ring system (RS) was done using RDKit. First, all monocyclic rings were retrieved; monocyclic rings were then fused depending on if individual ring systems shared atoms or not.

**Generative models.** The models studied in this study include CharRNN, REINVENT, AAE, VAE, ORGAN, LatentGAN, and GraphINVENT. The REINVENT code available at the github.com/undeadpixel/reinvent-randomized repo[35,41] was used; the CharRNN, AAE, VAE, and ORGAN codes available at the MOSES GitHub repository[10,42] were used; the LatentGAN code available at the github.com/Dierme/latent-gan repository[11,43] was used. Unlike original implement of LatentGAN, we retrained the embedded Deep Drug Coder (DDC) model with randomly selected 3M molecules from GDB-13 as the encoder and decoder component of LatentGAN. The DDC code available at the github.com/pcko1/Deep-Drug-Coder repository[44] was used; finally, the GraphINVENT code available at the github.com/MolecularAI/GraphINVENT repository [38,45,46] was used. All methods except for GraphINVENT are string-based generative models, whereas GraphINVENT is a graph-based generative model.

**Training.** The GDB-13 database is used as the reference chemical space for this study.[31] A one million (1M) molecule subset of GDB-13 was randomly selected and used as the training set for all the generative models. Another 200K molecules of GDB-13 were selected as the validation set for calculating the validation loss. During training, a check point model was saved at every epoch. The check point model with the lowest validation loss was chosen as the final model for sampling 1 billion (1B) SMILES.

**Hyperparameters**. For REINVENT, hyperparameters were taken from the GitHub repo.[41] For CharRNN, AAE, VAE, and ORGAN, the parameters were taken from the models' config file in MOSES GitHub repo without further optimization. For LatentGAN, , the default values of parameters in the GitHub repo were adopted.[11, 43] For GraphINVENT, parameters and hyperparameters for the best performing model (cGGNN) in the original publication were used and not further optimized.[46] Detailed hyperparameters for each model can be found in SI Table S2.

**Sampling**. Once each model was trained, 1B compounds were sampled from each trained model. The functional groups and ring systems were then identified for all sample sets as well as the full GDB-13 set. All compounds in the analysis were standardized by converting to RDKit canonical SMILES. Molecular graphs generated using GraphINVENT were further sanitized via canonicalizing and aromatizing during the conversion to canonical RDKit SMILES for a more fair comparison to the other models.
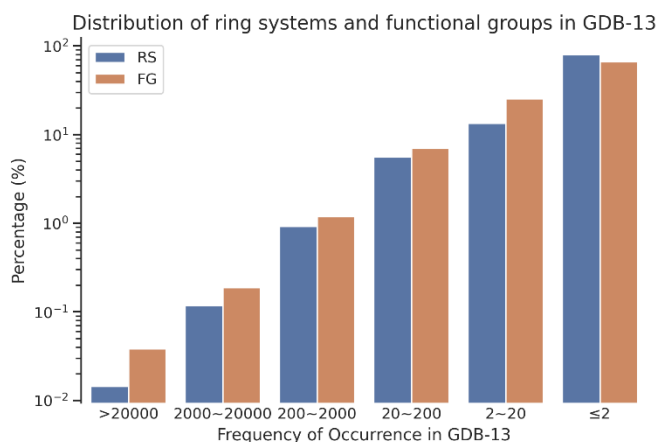
**Technical details.** For models in the MOSES repository and REINVENT, the training was done using Python 3.6[47] and PyTorch 1.4[48]. To accelerate sampling for 1B SMILES, the largest batch size allowed by the GPU memory was adopted; for example, ORGAN, AAE, and VAE adopted a sampling batch size of 25, 000, and CharRNN adopted a sampling batch size of 20,000. Also, LatentGAN was trained using tensorflow-gpu 2.2, which adopted a sampling batch size of 50,000.

153  All the computations were performed on Linux workstations with GeForce RTX 2080Ti graphic

154  cards using CUDA 10.1. Canonical SMILES and dataset analysis were carried out using RDKit.[39]

155  The Wasserstein distances[49] between distributions in Figure 2 were calculated with an in-house

156  script using SciPy.[50] Finally, GraphINVENT runs using Python 3.6 and PyTorch 1.2.

## Results and Discussions

### Analysis of the GDB-13 database

159  GDB-13 contains theoretically drug-like compounds whose heavy atom count is less than or equal

160  to 13 and, in total, comprises of 975,820,210 molecules, 21,852,845 ring systems, and 4,401,506

161  functional groups. The distribution of the occurrence frequency of these ring systems and

162  functional groups is shown in Figure 1. Figure 1 indicates that ~80% of ring systems and ~66%

163  functional groups in GDB-13 occur in compounds only 1-2 times, while only ~1% of ring systems

164  and functional groups are observed in GDB-13 molecules more than 200 times. In general, most

165  of the ring systems (~93%) and functional groups (~91%) appear in GDB-13 less than 20 times.



166

167  **Figure 1.** Distribution of ring systems (RS) and functional groups (FG) in GDB-13 according to

168  the frequency of occurrence. Y-axis is the percentage plotted on a logarithmic scale.
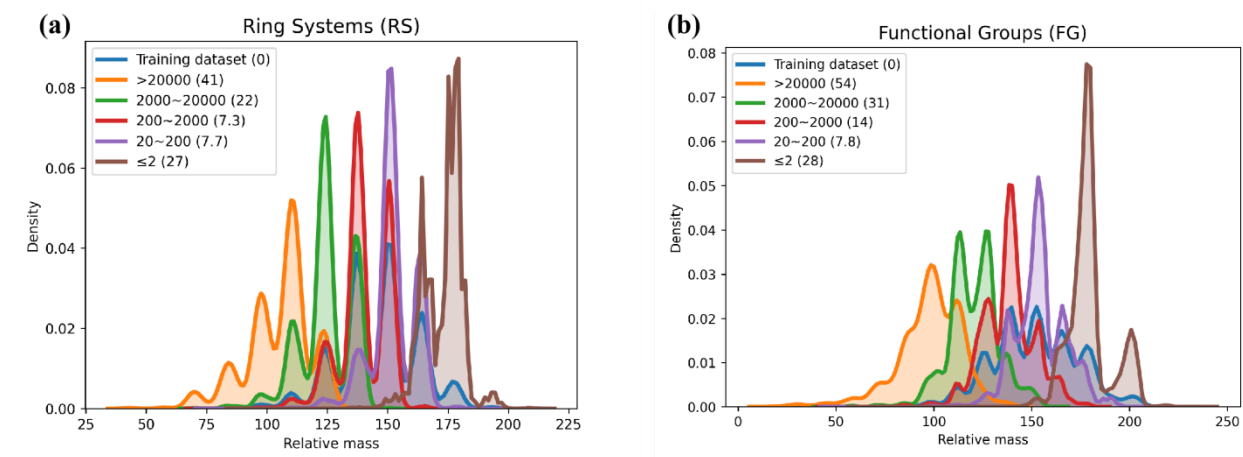
169    **Analysis of the 1M training dataset**

170    One million SMILES were randomly selected from the GDB-13 database for the training set,

171    which corresponds to roughly 0.1% of the total GDB-13 dataset. The training set contains around

172    0.9% of the ring systems and functional groups in the whole GDB-13 database (Table 1). The

173    coverage of the ring systems and functional groups is nine times as high as the coverage of

174    compounds, which is obviously due to the fact that some ring systems and functional groups occur

175    far more than once in GDB-13, as shown in Figure 1.

176    **Table 1.** Summary of GDB-13 coverage in the training set, consisting of 1M randomly selected

177    molecules.

| Item | Counts in the training dataset (1M) | Coverage of GDB-13 |
|---|---|---|
| Compounds | 1,000,000 | ~0.1% |
| Ring systems | 202,848 | ~0.9% |
| Functional groups | 38,209 | ~0.9% |

178

179



180    **Figure 2.** The frament weight distributions for the different substructures in GDB-13. The different

181    colors indicate distributions involving substructures that occur in GDB-13 a similar number of

182    times (i.e. orange is substructures that occur >20,000 times in GDB-13, brown is substructures that

183     occur ≤ 2 times). In the key, the numbers in parentheses indicate the Wasserstein distance between

184     the training set distribution and the indicated distribution. (a) Ring systems (RS). (b) Functional
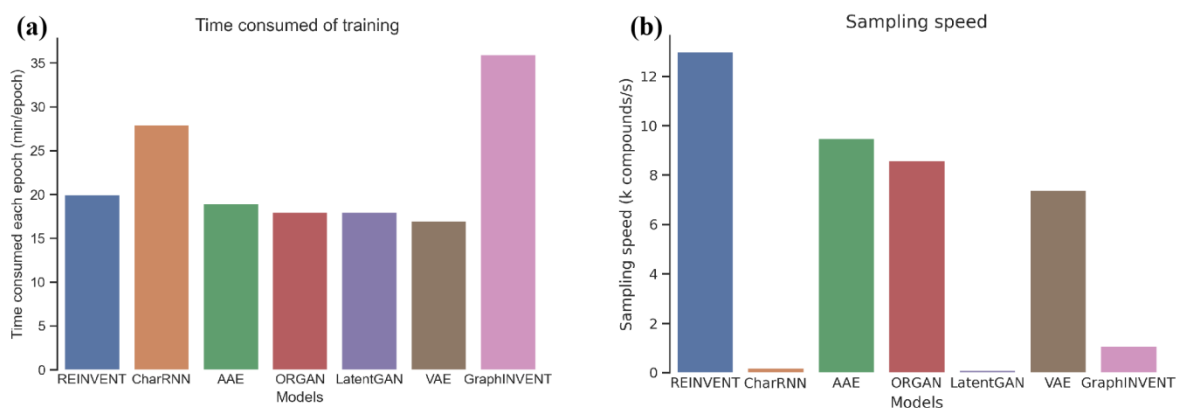
185     groups (FG).

186     The fragment weights (FWs) of ring systems and functional groups in the training set, grouped by

187     frequency of occurrence, are shown in Figure 2. The FWs here were calculated from the

188     composition of specific ring systems and functional groups rather than the full compound.  It is

189     observed that their probability of occurrence decreases with increasing FW. For example, the mean

190     FW of RS and FG which occur with a frequency >20,000 is around 100; however, for RS and FG

191     which occur <=2 times in GDB-13, the mean FW is around 170. More basic RS and FG, such as

192     C1CC1 (cyclopropyl) and C=O (carbonyl), respectively, tend to have smaller FW compared to

193     complex RS and FG. Furthermore, many complex RS and FG can be built from the basic

194     components via enumeration and combination following the chemical rules extracted from the

195     training dataset.

196     **Training and sampling speed**

197     All deep molecular generative models were trained with the same training set of 1M compounds.

198     Each epoch of training took 17-20 min for most models (Figure 3), except CharRNN (28 min) and

199     GraphINVENT (36 min). In general, the training speed of all the models is acceptable. We

200     observed that training SMILES-based models is faster than the graph-based model; this is

201     understandable because the action space of the graph-based model is much larger than any of the

202     SMILES-based models.

203     Nonetheless, the sampling speed of the generative models was observed to vary significantly. The

204     sampling speed of REINVENT, AAE, ORGAN, and VAE were all above 7000 compounds per

205     second, while the sampling speed of CharRNN, LatentGAN, and GraphINVENT were only 200,

206  100, and 1100 compounds per second, respectively. Notably, CharRNN and REINVENT share

207  similar architecture of character-level recurrent neural networks. The difference of their

208  performance is mainly due to CharRNN implementation provided by MOSES adopts a larger size

209  of architecture. The detailed hyper parameters are given as Table S2 in the supporting materials.

210  It should also be noted that both training and sampling speeds are strongly related to the batch size

211  that is limited by the memory of the GPU. In current work, the default batch size as specified in

212  the code was used during the training, while for sampling, the largest batch size allowed by the

213  GPU memory was chosen.

214  Given the relatively small size of the training set (1M molecules), all the deep generative models

215  had a tractable training speed. In terms of sampling, the sampling speed was limited by each

216  model's architecture and size; using a larger sampling batch size allowed by a greater GPU

217  memory could boost the sampling speed.



218
219  **Figure 3.** Training and sampling speeds of the generative models benchmarked in this work. (a)

220  Time consumed per epoch during training. (b) Sampling speed, which is the number of

221  SMILES/graphs generated per second (including invalid ones).

222  **Validity and repetition rate of sampled molecules**

223     We first check the validity of the molecules generated by all the deep generative models, which is

224     defined as the percentage of chemically valid SMILES/graphs in the 1B generated set. Table 2

225     shows that the validity in general is satisfactory for all models, where most models achieve a

226     validity higher than 90 percent. RNN based models (REINVENT and CharRNN) have the highest

227     validity which is above 99.3% (Table 2). The validity of LatentGAN and GraphINVENT are 85.4%

228     and 95.3% respectively, which are lowest among all the models. In order to check how much

229     duplication is generated among the sample sets, the repetition rate ($R_{rept}$) was calculated via the

230     formula below:

231 
$$R_{rept} = \frac{N_{valid} - N_{unique}}{N_{unique}}, \tag{1}$$

232     where, $N_{valid}$ is the number of total valid molecules in the 1B generated set and $N_{unique}$ is the

233     number of unique valid molecules in the 1B generated set (i.e. duplicates removed). The compound

234     repetition rates of most deep generative models were around 1.0, that is to say, most compounds

235     were sampled twice on the average. ORGAN and CharRNN have the highest repetition rates,

236     which are 3.8 and 1.4 respectively, whereas GraphINVENT and LatentGAN have the lowest (0.7).

237     It seems that all the deep generative models had a satisfactory high percent validity that was above

238     85% in this study. The validity of CharRNN reached as high as 99.7%. ORGAN had a repetitive

239     rate as high as 3.8, which means that each generated compound was sampled 4.8 times on average.

240     The high repetition rate resulted in a low overall compound coverage for ORGAN, where the

241     coverage was as low as 16%.

242     **Table 2.** Percentage of the valid molecules and molecular repetition rate in the 1B generated set

243     for each model in this study. The uncertainty in the percent validity was less than a fraction of a

244     percentage point for each model.

| Model | REINVENT | CharRNN | AAE | ORGAN | LatentGAN | VAE | GraphINVENT |
|---|---|---|---|---|---|---|---|
| Validity (%) | 99.3 | 99.7 | 97.8 | 97.2 | 85.4 | 98.2 | 95.3 |
| Repetition rate | 0.9 | 1.4 | 0.9 | 3.8 | 0.7 | 1.0 | 0.7 |

245

246 **Coverage of GDB-13 chemical space**

247 The molecule and substructure coverage of GDB-13 space for all generative models studied herein

248 is shown in Figure 4a. It can be seen that all the models possess good capabilities for generalization,

249 surpassing the coverage of the 1M training set used, which has a ~0.1% coverage of GDB-13

250 compounds, ~0.9% coverage of GDB-13 ring systems, and ~0.9% coverage of GDB-13 functional

251 groups. REINVENT achieves the highest compound and FG coverage (39% and 26%,

252 respectively), while AAE achieves best RS coverage (41%). The GAN models (ORGAN and

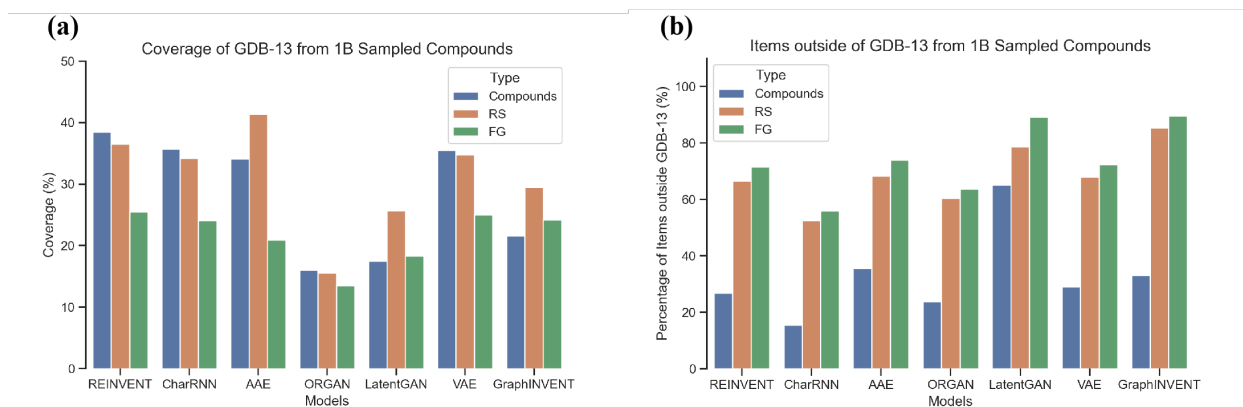253 LatentGAN) have lowest coverage at all three levels.

254 Using these new metrics, the difference in performance among these models is more pronounced;

255 this is in contrast to a previous benchmarking study using the MOSES metrics,[10] where the two

256 GAN models appear to perform similarly with the CharRNN, AAE, and VAE models.

257 Overall, REINVENT, CharRNN, AAE, and VAE are the top-ranking models in this benchmarking

258 study. They have a compound coverage, RS coverage, and FG coverage around 34%, 34%, and

259 21%, respectively, in all cases. The performance of GraphINVENT is in the middle rank among

260 the generative models in this study, and demonstrates coverage scores of 22%, 30%, and 24% for

261 compound coverage, RS coverage, and FG coverage, respectively.

262

263

264

265



**Figure 4.** Coverage of GDB-13 chemical space using 1B sampled molecules. (a) Coverage of compounds, ring systems (RS), and functional groups (FG) in GDB-13 ($P_{covered}$). (b) Percentage of sampled molecules, RS, and FG that are outside the chemical space of GDB-13 ($P_{out}$).

Coverage of compounds, RS, and FG in GDB-13 was calculated via the formula below:

$$P_{covered} = \frac{N_{unique\_in}}{N_{GDB13}} * 100\%, \tag{2}$$

where $N_{unique\_in}$ is the number of unique valid sampled compounds, RS, or FG that are also found in GDB-13, and $N_{GDB13}$ is the total number of compounds, RS, or FG present in GDB-13.

The percentage of sampled compounds, RS, or FG that are outside the chemical space of GDB-13 was calculated via the formula below:

$$P_{out} = \frac{N_{unique\_out}}{N_{unique}} * 100\%, \tag{3}$$

where $N_{unique\_out}$ is the number of unique valid sampled compounds, RS, or FG that are *not* found in GDB-13, and $N_{unique}$ is the total number of unique valid compounds, RS, or FG in the generated sets.

14

280     There are four major metrics mentioned above, namely validity, repetition rate, coverage of GDB-

281     13 chemical space, and percentage outside GDB-13. Validity represents how good a generative

282     model has learned the chemical rules for constructing compounds; repetition rate represents how

283     much structure duplication exists in the generated compound set; generalization capacities of

284     models can be measured with the coverage of GDB-13 after being trained on a smaller fraction of

285     chemical space. As a supplement to above metrics, percentage outside GDB-13 shows how many

286     sampled compounds fall outside the scope of GDB-13 (which are usually non drug-like

287     compounds). Also, these four metrics are not independent from each other. For example, if a model

288     has a high validity and a small percentage sampled outside GDB-13, given that exactly 1B

289     compounds are sampled, the only reasonable explanation for a low GDB-13 coverage is a high

290     repetition rate.

291     Figure 4b shows the generated structures outside GDB-13. As GDB-13 uses filters to remove

292     molecules that do not satisfy simple chemical stability and synthetic feasibility rules, such as ring-

293     strain criteria and valency rules, there are many structures that can be generated which violate the

294     filters used by GDB-13. For example, there are around 27% valid SMILES generated by

295     REINVENT which fall outside the scope of GDB-13 chemical space. However, for CharRNN,

296     only 15% of its respective generated sets fall outside GDB-13, which is lower than other models

297     in this study. As the percent validity of the structures generated by both models is above 97%, we

298     conclude that the lower fraction of compounds outside GDB-13 is due to the high repetition rate

299     of compounds for these models, as shown in Table 2. As for the percentage of RS and FG outside

300     of the scope of GDB-13, more than 50% of all FG and RS found in the generated sets for each

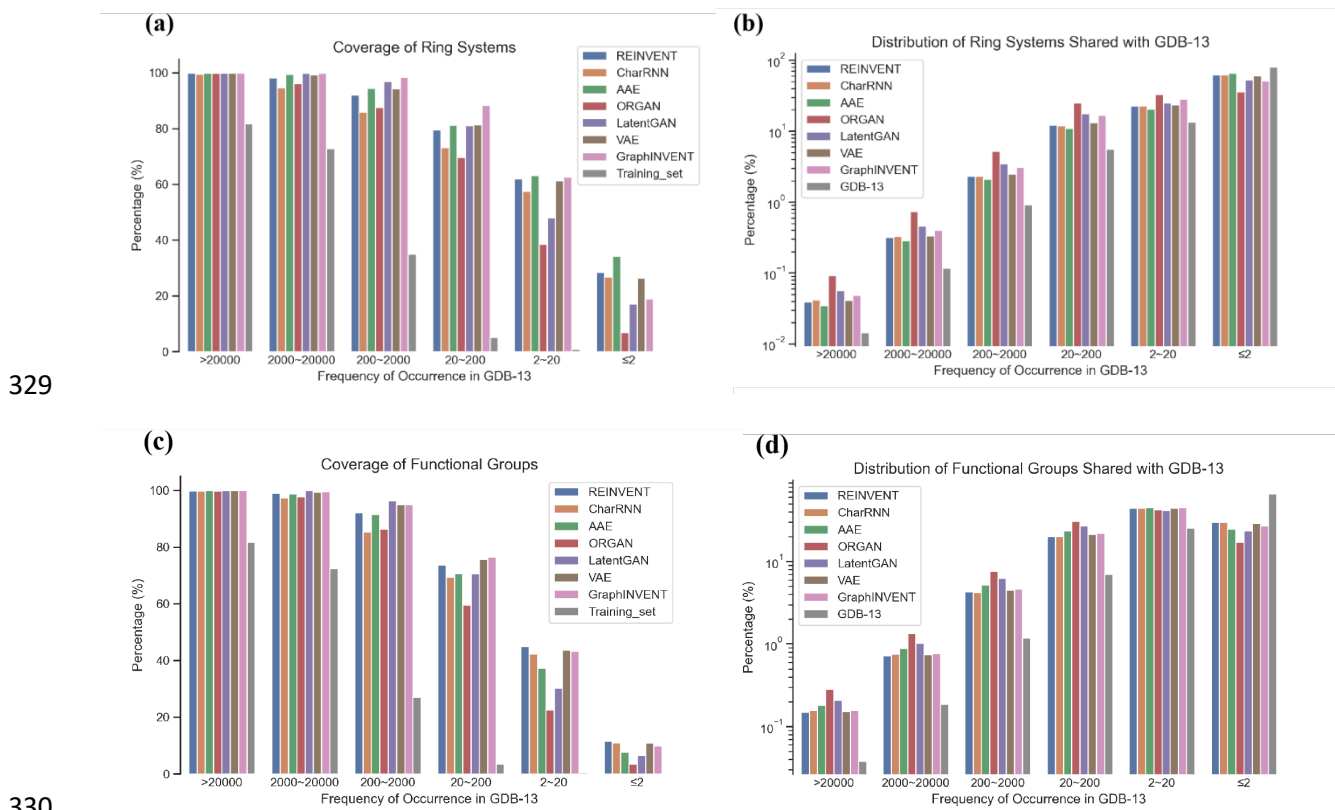301     model are outside the GDB-13 chemical space.

302     After training with a subset of the GDB-13 database (0.1%), all the generative models showed

303     promising performance in terms of compound coverage. Around 16% compounds of GDB-13 were

304     covered with 1B SMILES sampled by the LatentGAN, which is 160 times greater than the

305     coverage of the training dataset itself. The model with the best performance in this study is

306     REINVENT, which has an observed compound coverage as high as 39%. Thus, we conclude that

307     deep generative models in general have satisfactory learning and generalization capacities. In

308     terms of overall GDB-13 compound coverage, the rank of performance in descending order is

309     REINVENT > CharRNN > VAE > AAE > GraphINVENT > LatentGAN > ORGAN.

310     The GDB-13 coverage of RS and FG was generally less than the coverage of compounds, except

311     in the cases of AAE, LatentGAN, and GraphINVENT. However, in these cases, greater than 60%

312     RS and FG in the generated set were outside the scope GDB-13 chemical space, while less than

313     40% of generated molecules were outside GDB-13 (except LatentGAN). In terms of RS coverage

314     of GDB-13, the rank of performance in descending order is AAE > REINVENT > VAE >

315     CharRNN > GraphINVENT > LatentGAN > ORGAN. In terms of FG coverage of GDB-13, the

316     rank of performance in descending order is REINVENT > VAE > GraphINVENT > CharRNN >

317     AAE > LatentGAN > ORGAN. Examples of the most commonly observed groups in structures

318     generated by the two best models in terms of functional groups and ring systems recovery,

319     REINVENT and AAE, are shown in Figures 6 & 8. Examples of the most commonly observed

320     groups that are outside of GDB-13 in structures and generated by LatentGAN, are shown in Figures

321     7 & 9.

322     It is worthwhile to mention that the original LatentGAN adopts a heteroencoder and decoder model

323     (DDC) trained on ChEMBL dataset, the LatentGAN had a compounds coverage, RS coverage and

324     FG coverage of GDB-13 as 13%, 15% and 18%, respectively. When the DDC model were trained

325    on a 3M subset of GDB-13 instead, the compounds coverage, RS coverage and FG coverage of

326    GDB-13 increased to 18%, 26% and 18%, respectively. Thus, we adopted the heteroencoder and

327    decoder model trained on the 3M subset in this study.

**Relationship between the coverage of GDB-13 and occurrence frequency**

329



330

331    **Figure 5.** Coverage of GDB-13 chemical space from 1B sampled molecules, grouped by the

332    occurrence frequency of molecules in GDB-13. (a & c) Coverage of RS and FG. (b & d)

333    Distribution of generated RS and FG that are shared with the chemical space of GDB-13. The y-

334    axes for (b) and (d) are displayed in logarithmic scale.

335    The coverage of GDB-13 chemical space from 1B sampled molecules, grouped by the

336    occurrence frequency of molecules in GDB-13, was calculated via the formula below:

17

$$P_{covered} = \frac{N_{unique\_in}(R_m, \ R_n)}{N_{GDB13}(R_m, \ R_n)} * 100\%, \tag{4}$$

337

338 where $N_{unique\_in}(R_m, \ R_n)$ is the number of unique RS or FG in the sampled set that have an

339 occurrence frequency in the interval of $R_m - R_n$ (including $R_n$) in GDB-13, and $N_{GDB13}(R_m, \ R_n)$

340 is the total number of RS or FG in GDB-13 with an occurrence frequency in the interval of $R_m -$

341 $R_n$ (including $R_n$). As such, $P_{covered}$ represents the coverage of specific set of substructures

342 $N_{GDB13}(R_m, \ R_n)$ of GDB-13 from the 1B generated set.

343 In Figures 5a and 5c, the RS and FG coverage of various models is broken down into different

344 frequency sections to examine the coverage performance for different types of substructures.

345 Figure 5 shows that for high frequency RS and FG, the coverage is high and quite similar among

346 all models, while for less frequent RS and FG, the coverage reveals differences between models.

347 On the other hand, comparing with the training set, all models demonstrate clear enrichment of RS

348 and FG coverage, and the enrichment gets bigger as the RS and FG frequency is lower. As for RS

349 and FG at the occurrence ranges of ">20000", "2000-20000", and "200-2000", the coverage is

350 close to 100% for all models, while the coverage of the training dataset is around 82%, 73%, and

351 31% at these respective occurrence frequency ranges. As for RS at the occurrence range of "20-

352 200", "2-20" and "≤2", most generative models have an RS coverage of around 80%, 60%, and

353 30%, compared to only 5%, 1%, and 0% for the training dataset. The coverage of FG at the

354 different occurrence frequency ranges has a similar pattern to the RS coverage.

355 Similarly, distribution of generated RS and FG that are shared with the chemical space of GDB-

356 13 was calculated via the formula below:

$$P_{dist} = \frac{N_{unique\_in}(R_m, \ R_n)}{N_{unique\_in}} * 100\%, \tag{5}$$

357

358    where $N_{unique\_in}(R_m, R_n)$ is the number of unique RS or FG in the sampled set that have an

359    occurrence frequency in the range of $R_m$ to $R_n$ in GDB-13, and $N_{unique\_in}$ is the total number of

360    unique RS or FG in the generated set, which are also included in GDB-13. Thus, $P_{dist}$ is a metric

361    of the distribution of RS or FG that are shared with GDB-13 at different occurrence ranges.

362    The distributions of generated RS and FG corresponding to occurrence frequency in GDB-13 are

363    shown in Figures 5b & 5d. Given that most RS and FG have an occurrence frequency below 20 in

364    the GDB-13 database (as shown in Figure 1), the overall coverage of RS and FG is thus dominated

365    by ones with low occurrence frequency.

366    The most frequent and least frequent ring systems and functional groups sampled by the deep

367    generative models are listed in Figures 6-9. The most often sampled ring systems are simple carbon

368    cycles or aromatic heterocycles containing O and N atoms, such as C1CC1 (cyclopropane), which

369    was sampled up to 78M times in the 1B sample set, and C1COC1 (oxetane), which were sampled

370    up to 26M times in the 1B sample set. For comparison, the benzene ring ranked 85[th] among the

371    most common sampled ring systems. As for the least common sampled ring systems, they were

372    usually complex macrocycles that were only sampled once out of the 1B compounds generated.

373    The most commonly sampled functional groups are ordinary small ones, such as single oxygen

374    and nitrogen atoms, C-C double bonds, and C-C triple bonds. The least commonly sampled

375    functional groups are those with complex structures formed by a combination of simple ones. The

376    ring systems and functional groups that are not included in GDB-13 usually do not conform to

377    simple chemical stability and synthetic feasibility rules.

378    Most of the RS (~93%) and FG (~91%) found in the generated sets that are also found in GDB-13

379    are seen less than 20 times. As the results show in Figure 2, RS and FG that occur more frequently

380    in GDB-13 tend to have smaller fragment weights. The building blocks of RS and FG are basic

381    rings and functional groups with simple structures and small fragment weights. More complex RS

382    and FG can be built via the combination of these basic components.
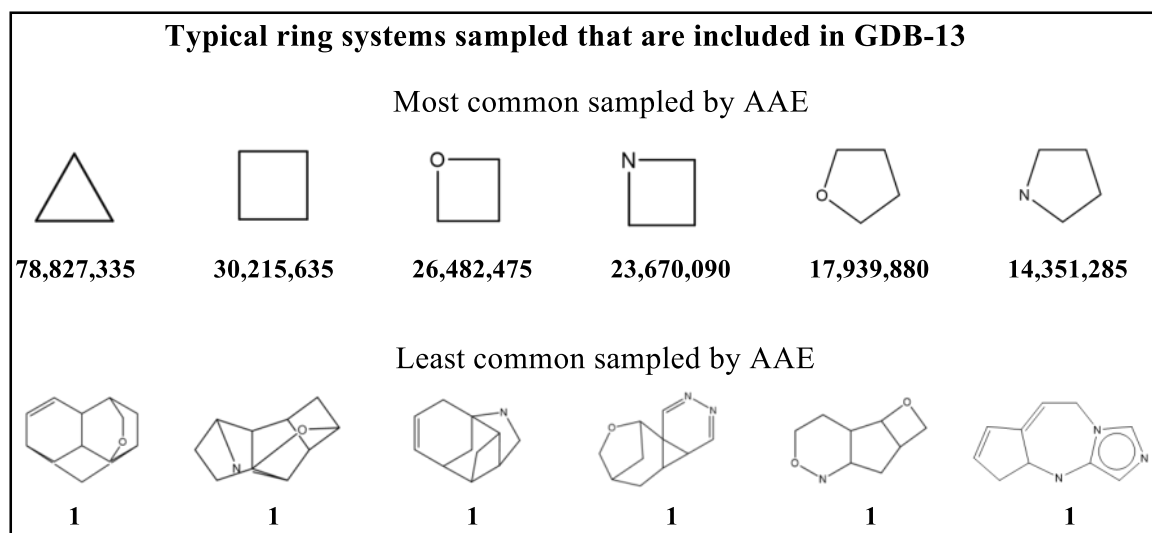
383    The coverage of RS and FG with an occurrence frequency in GDB-13 greater than 200 was nearly

384    100%. This is because these RS and FG can be easily obtained via combinations of smaller

385    fragments. However, given that as many as up to 13 heavy atoms were considered in constructing

386    the GDB-13 database, most RS and FG possess complex structures and were included in

387    compounds of GDB-13 less than 20 times. RS and FG that occur less than 20 times in the generated

388    sets dominate the coverage of the deep generative models.

389    Besides, as shown in Figures 10, most common ring systems and functional groups sampled by

390    generative models have close relative occurrence frequency compared to their distribution in
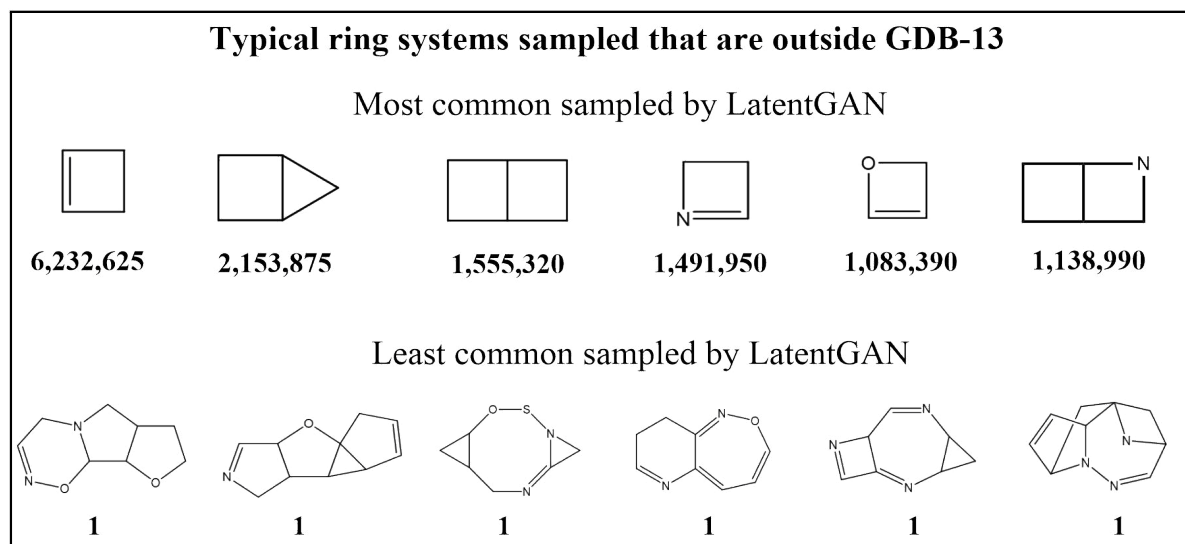
391    GDB-13.

392    **Model comparison**

393    It is interesting to observe that these models describe the chemical space so differently, although

394    trained with the same training set. It seems that the RS and FG coverage of GraphINVENT is

395    higher than its overall molecular coverage, one reason could be due to its large action space; that

396    is, the number of possible "correct" sampled actions at any stage during graph generation is much

397    larger than it is for SMILES-based methods which must use only tokens sampled in the training

398    set. As such, given that GraphINVENT samples actions probabilistically, it is possible that

399    sequences of actions are sampled which have never been seen in the training set, thus leading to

400    new molecules. Another interesting observation is that GAN based models generally perform

401    worst in terms of GDB-13 coverage on all three metrics, one reason could be due to that, in the

402    adversarial training, the generator is supposed to mimic the true data as much as possible to fool

403    the discriminator, which deteriorates its generalization capability to a certain extent. We also
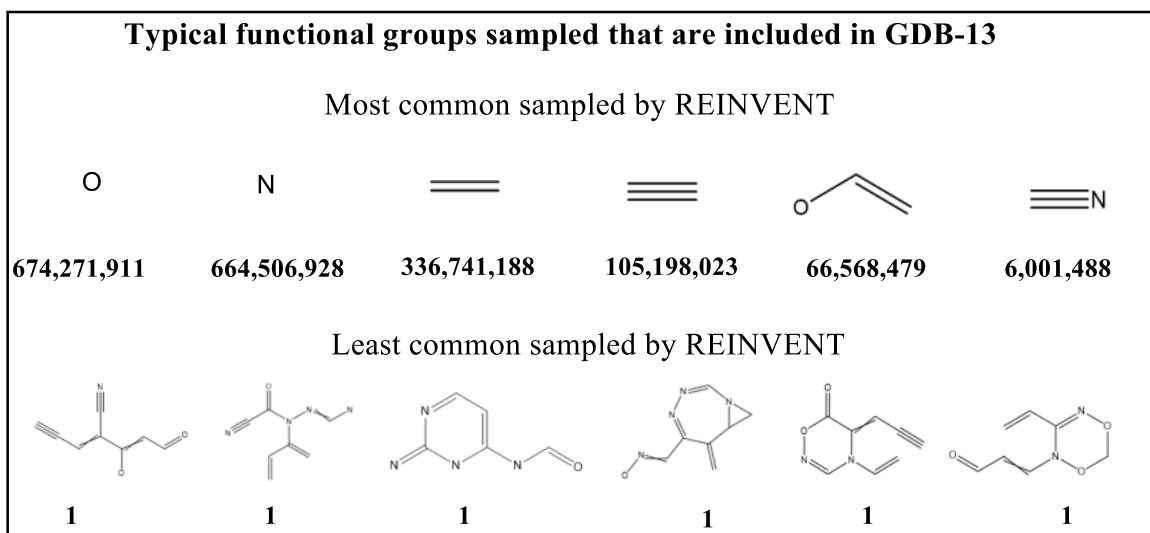
404     noticed that the performance of REINVENT and CharRNN is somehow similar, while their

405     sampling speed has very large difference. Given that both models are based on the same RNN

406     architecture, suggesting that the technical implementation of CharRNN is suboptimal.



**Typical ring systems sampled that are included in GDB-13**

Most common sampled by AAE

| 78,827,335 | 30,215,635 | 26,482,475 | 23,670,090 | 17,939,880 | 14,351,285 |

Least common sampled by AAE

| 1 | 1 | 1 | 1 | 1 | 1 |

407

408     **Figure 6.** Typical ring systems that are sampled by AAE, which are included in GDB-13. The

409     numbers below the structures in the figure are the occurrence frequency of ring systems in the 1B

410     sampled compounds.



**Typical ring systems sampled that are outside GDB-13**

Most common sampled by LatentGAN

| 6,232,625 | 2,153,875 | 1,555,320 | 1,491,950 | 1,083,390 | 1,138,990 |

Least common sampled by LatentGAN

| 1 | 1 | 1 | 1 | 1 | 1 |

411

412     **Figure 7.** Typical ring systems that are sampled by LatentGAN, which are outside GDB-13.

**Typical functional groups sampled that are included in GDB-13**

Most common sampled by REINVENT

674,271,911     664,506,928     336,741,188     105,198,023     66,568,479     6,001,488

Least common sampled by REINVENT

1     1     1     1     1     1
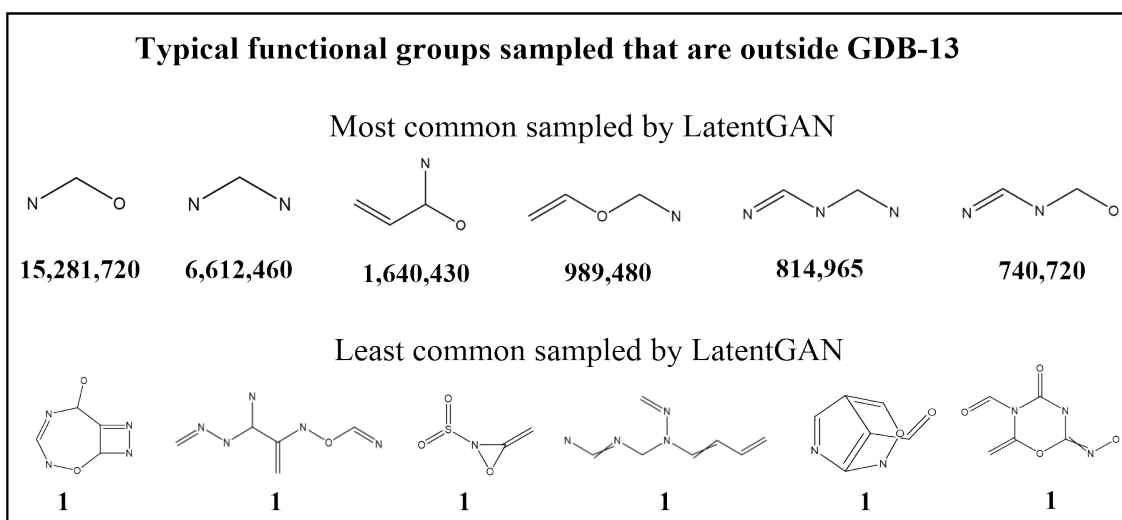
**Figure 8.** Typical functional groups that are sampled by REINVENT, which are included in GDB-13.



**Typical functional groups sampled that are outside GDB-13**

Most common sampled by LatentGAN

15,281,720     6,612,460     1,640,430     989,480     814,965     740,720

Least common sampled by LatentGAN

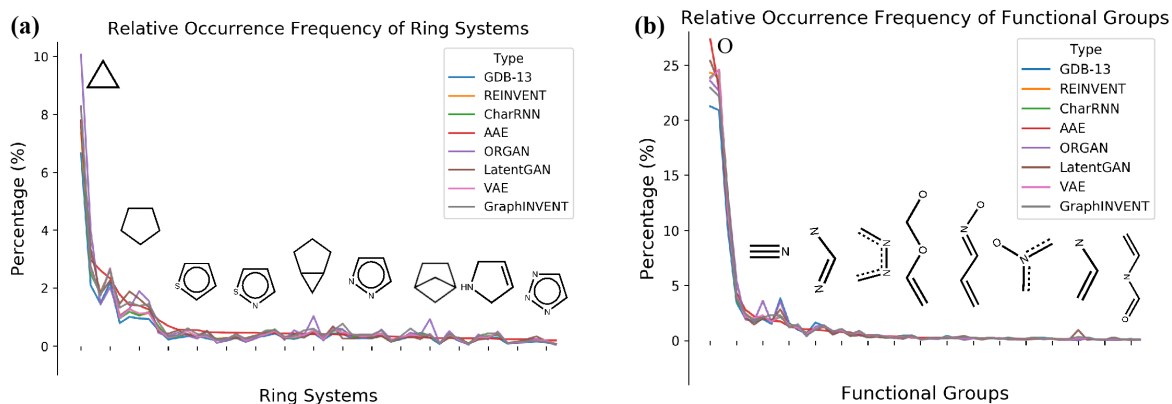1     1     1     1     1     1

**Figure 9.** Typical functional groups that are sampled by LatentGAN, which are outside GDB-13.

**(a)** Relative Occurrence Frequency of Ring Systems

**(b)** Relative Occurrence Frequency of Functional Groups

419

**Figure 10.** Relative occurrence frequency of most common functional groups and ring systems.

## Conclusions

Molecules consist of a variety of ring systems and functional groups, which are connected in different ways to form molecules. The most basic ring systems and functional groups have simple structures and small fragment weights; these can be found in GDB-13 molecules over dozens of times. More complex ring systems and functional groups have complicated structures and large fragment weights, and might only occur in GDB-13 a handful of times. However, due to their structural variety and enormous quantity (>90%), complex ring systems and functional groups are strong components affecting the coverage of GDB-13.

All the deep generative models studied in this work have over 100 times greater chemical space coverage for GDB-13 using 1B samples than the training set (1M) used to train the models. In terms of compound coverage of GDB-13, the best model (REINVENT) reached ~39% coverage, far beyond the coverage of ORGAN (~16%), which ranked lowest amongst the models in this study. Depending on the generative task, the deep generative model used should thus be chosen

434 carefully, as there are differences in how all these seemingly similar models sample the chemical

435 space.

## Associated Content

## Author Information

### Corresponding Author

439 Hongming Chen, E-mail: chen_hongming@grmh-gdl.cn

### Author Contributions

441 J. Z. ran training and generation jobs using REINVENT, CharRNN, VAE; LatentGAN, and

442 ORGAN. R. M. ran training and generation jobs using GraphINVENT. R. M. and J. Z. ran

443 benchmarking calculations for this work, and J. Z. made all figures. The manuscript was written

444 through the contributions of all authors. All authors have given approval to the final version of the

445 manuscript.

### Supplementary Materials

450 The detailed hyperparameters and training loss curves of all models can be found in supplementary

451 materials. The training, sampling, and analysis script could found in the GitHub repository,

452 https://github.com/jeah-z/Generative_Models_benchmark_gdb13.

## Abbreviations

454    RS, Ring system(s); FG, Functional group(s); GAN, Generative adversarial network; GNN,

455    graph neural network; RNN, recurrent neural network.

## References

457    (1) Ciregan, D.; Meier, U.; Schmidhuber, J. Multi-Column Deep Neural Networks for Image Classification.
458    In 2012 IEEE conference on computer vision and pattern recognition, 2012; IEEE: 2012; pp 3642-3649.
459    (2) Krizhevsky, A.; Sutskever, I.; Hinton, G. E. Imagenet Classification with Deep Convolutional Neural
460    Networks. In Advances in neural information processing systems, 2012; 2012; pp 1097-1105.
461    (3) Taigman, Y.; Yang, M.; Ranzato, M. A.; Wolf, L. Deepface: Closing the Gap to Human-Level Performance
462    in Face Verification. In Proceedings of the IEEE conference on computer vision and pattern recognition,
463    2014; 2014; pp 1701-1708.
464    (4) Silver, D.; Huang, A.; Maddison, C. J.; Guez, A.; Sifre, L.; Van Den Driessche, G.; Schrittwieser, J.;
465    Antonoglou, I.; Panneershelvam, V.; Lanctot, M., Mastering the Game of Go with Deep Neural Networks
466    and Tree Search. *nature* **2016**, *529*, 484-489.
467    (5) Hadjeres, G.; Pachet, F.; Nielsen, F. Deepbach: A Steerable Model for Bach Chorales Generation. In
468    International Conference on Machine Learning, 2017; PMLR: 2017; pp 1362-1371.
469    (6) Garg, S.; Rish, I.; Cecchi, G.; Lozano, A., Neurogenesis-Inspired Dictionary Learning: Online Model
470    Adaption in a Changing World. *arXiv preprint arXiv:1701.06106* **2017**.
471    (7) Johnson, M.; Schuster, M.; Le, Q. V.; Krikun, M.; Wu, Y.; Chen, Z.; Thorat, N.; Viégas, F.; Wattenberg,
472    M.; Corrado, G., Google's Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation.
473    *Transactions of the Association for Computational Linguistics* **2017**, *5*, 339-351.
474    (8) Bjerrum, E. J.; Threlfall, R., Molecular Generation with Recurrent Neural Networks (Rnns). *arXiv*
475    *preprint arXiv:1705.04612* **2017**.
476    (9) Kotsias, P.-C.; Arús-Pous, J.; Chen, H.; Engkvist, O.; Tyrchan, C.; Bjerrum, E. J., Direct Steering of De
477    Novo Molecular Generation with Descriptor Conditional Recurrent Neural Networks. *Nature Machine*
478    *Intelligence* **2020**, *2*, 254-265.
479    (10) Polykovskiy, D.; Zhebrak, A.; Sanchez-Lengeling, B.; Golovanov, S.; Tatanov, O.; Belyaev, S.; Kurbanov,
480    R.; Artamonov, A.; Aladinskiy, V.; Veselov, M.; Kadurin, A.; Johansson, S.; Chen, H.; Nikolenko, S.; Aspuru-
481    Guzik, A.; Zhavoronkov, A., Molecular Sets (Moses): A Benchmarking Platform for Molecular Generation
482    Models. *Front Pharmacol* **2020**, *11*, 565644.
483    (11) Prykhodko, O.; Johansson, S. V.; Kotsias, P.-C.; Arús-Pous, J.; Bjerrum, E. J.; Engkvist, O.; Chen, H., A
484    De Novo Molecular Generation Method Using Latent Vector Based Generative Adversarial Network.
485    *Journal of Cheminformatics* **2019**, *11*, 74.
486    (12) Schneider, G.; Fechner, U., Computer-Based De Novo Design of Drug-Like Molecules. *Nat Rev Drug*
487    *Discov* **2005**, *4*, 649-63.
488    (13) DiMasi, J. A.; Grabowski, H. G.; Hansen, R. W., Innovation in the Pharmaceutical Industry: New
489    Estimates of R&D Costs. *Journal of health economics* **2016**, *47*, 20-33.
490    (14) Patel, H.; Bodkin, M. J.; Chen, B.; Gillet, V. J., Knowledge-Based Approach to De Novo Design Using
491    Reaction Vectors. *Journal of chemical information and modeling* **2009**, *49*, 1163-1184.
492    (15) Schneider, G.; Lee, M.-L.; Stahl, M.; Schneider, P., De Novo Design of Molecular Architectures by
493    Evolutionary Assembly of Drug-Derived Building Blocks. *Journal of computer-aided molecular design* **2000**,
494    *14*, 487-494.
495    (16) Spiegel, J. O.; Durrant, J. D., Autogrow4: An Open-Source Genetic Algorithm for De Novo Drug Design
496    and Lead Optimization. *Journal of Cheminformatics* **2020**, *12*, 1-16.

497 (17) Chen, H.; Engkvist, O.; Wang, Y.; Olivecrona, M.; Blaschke, T., The Rise of Deep Learning in Drug
498 Discovery. *Drug discovery today* **2018**, *23*, 1241-1250.
499 (18) Blaschke, T.; Arus-Pous, J.; Chen, H.; Margreitter, C.; Tyrchan, C.; Engkvist, O.; Papadopoulos, K.;
500 Patronov, A., Reinvent 2.0: An Ai Tool for De Novo Drug Design. *J Chem Inf Model* **2020**, *60*, 5918-5922.
501 (19) Blaschke, T.; Olivecrona, M.; Engkvist, O.; Bajorath, J.; Chen, H., Application of Generative
502 Autoencoder in De Novo Molecular Design. *Molecular informatics* **2018**, *37*, 1700123.
503 (20) Guimaraes, G. L.; Sanchez-Lengeling, B.; Outeiral, C.; Farias, P. L. C.; Aspuru-Guzik, A., Objective-
504 Reinforced Generative Adversarial Networks (Organ) for Sequence Generation Models. *arXiv preprint*
505 *arXiv:1705.10843* **2017**.
506 (21) Kadurin, A.; Aliper, A.; Kazennov, A.; Mamoshina, P.; Vanhaelen, Q.; Khrabrov, K.; Zhavoronkov, A.,
507 The Cornucopia of Meaningful Leads: Applying Deep Adversarial Autoencoders for New Molecule
508 Development in Oncology. *Oncotarget* **2017**, *8*, 10883.
509 (22) Grisoni, F.; Moret, M.; Lingwood, R.; Schneider, G., Bidirectional Molecule Generation with Recurrent
510 Neural Networks. *Journal of chemical information and modeling* **2020**, *60*, 1175-1183.
511 (23) Gupta, A.; Müller, A. T.; Huisman, B. J.; Fuchs, J. A.; Schneider, P.; Schneider, G., Generative Recurrent
512 Networks for De Novo Drug Design. *Molecular informatics* **2018**, *37*, 1700111.
513 (24) Olivecrona, M.; Blaschke, T.; Engkvist, O.; Chen, H., Molecular De-Novo Design through Deep
514 Reinforcement Learning. *J Cheminform* **2017**, *9*, 48.
515 (25) Segler, M. H.; Kogej, T.; Tyrchan, C.; Waller, M. P., Generating Focused Molecule Libraries for Drug
516 Discovery with Recurrent Neural Networks. *ACS central science* **2018**, *4*, 120-131.
517 (26) Yuan, Q.; Santana-Bonilla, A.; Zwijnenburg, M. A.; Jelfs, K. E., Molecular Generation Targeting Desired
518 Electronic Properties Via Deep Generative Models. *Nanoscale* **2020**, *12*, 6744-6758.
519 (27) Kalchbrenner, N.; Grefenstette, E.; Blunsom, P., A Convolutional Neural Network for Modelling
520 Sentences. *arXiv preprint arXiv:1404.2188* **2014**.
521 (28) Kingma, D. P.; Welling, M., Auto-Encoding Variational Bayes. *arXiv preprint arXiv:1312.6114* **2013**.
522 (29) Kusner, M. J.; Paige, B.; Hernández-Lobato, J. M., Grammar Variational Autoencoder. *arXiv preprint*
523 *arXiv:1703.01925* **2017**.
524 (30) Brown, N.; Fiscato, M.; Segler, M. H.; Vaucher, A. C., Guacamol: Benchmarking Models for De Novo
525 Molecular Design. *Journal of chemical information and modeling* **2019**, *59*, 1096-1108.
526 (31) Blum, L. C.; Reymond, J.-L., 970 Million Druglike Small Molecules for Virtual Screening in the Chemical
527 Universe Database Gdb-13. *Journal of the American Chemical Society* **2009**, *131*, 8732-8733.
528 (32) Arus-Pous, J.; Blaschke, T.; Ulander, S.; Reymond, J. L.; Chen, H.; Engkvist, O., Exploring the Gdb-13
529 Chemical Space Using Deep Generative Models. *J Cheminform* **2019**, *11*, 20.
530 (33) Ertl, P.; Altmann, E.; McKenna, J. M., The Most Common Functional Groups in Bioactive Molecules
531 and How Their Popularity Has Evolved over Time. *Journal of Medicinal Chemistry* **2020**, *63*, 8408-8418.
532 (34) Taylor, R. D.; MacCoss, M.; Lawson, A. D., Rings in Drugs. *J Med Chem* **2014**, *57*, 5845-59.
533 (35) Arús-Pous, J.; Johansson, S. V.; Prykhodko, O.; Bjerrum, E. J.; Tyrchan, C.; Reymond, J.-L.; Chen, H.;
534 Engkvist, O., Randomized Smiles Strings Improve the Quality of Molecular Generative Models. *Journal of*
535 *Cheminformatics* **2019**, *11*.
536 (36) Gómez-Bombarelli, R.; Wei, J. N.; Duvenaud, D.; Hernández-Lobato, J. M.; Sánchez-Lengeling, B.;
537 Sheberla, D.; Aguilera-Iparraguirre, J.; Hirzel, T. D.; Adams, R. P.; Aspuru-Guzik, A., Automatic Chemical
538 Design Using a Data-Driven Continuous Representation of Molecules. *ACS central science* **2018**, *4*, 268-
539 276.
540 (37) Makhzani, A.; Shlens, J.; Jaitly, N.; Goodfellow, I.; Frey, B., Adversarial Autoencoders. *arXiv preprint*
541 *arXiv:1511.05644* **2015**.
542 (38) Mercado, R.; Rastemo, T.; Lindelöf, E.; Klambauer, G.; Engkvist, O.; Chen, H.; Bjerrum, E. J., Practical
543 Notes on Building Molecular Graph Generative Models. *ChemRxiv* **2020**, *Preprint*.
544 (39) Landrum, G., Rdkit: Open-Source Cheminformatics. **2006**.

545    (40) Ertl, P., An Algorithm to Identify Functional Groups in Organic Molecules. *Journal of cheminformatics*
546    **2017**, *9*, 1-7.
547    (41) Arús-Pous, J. Reinvent-Randomized. h[ttps://github.com/undeadpixel/reinvent-randomized](https://github.com/undeadpixel/reinvent-randomized)
548    (accessed Sep 1, 2020).
549    (42) Polykovskiy, D.; Zhebrak, A.; Sanchez-Lengeling, B.; Golovanov, S.; Tatanov, O.; Belyaev, S.; Kurbanov,
550    R.; Artamonov, A.; Aladinskiy, V.; Veselov, M. Moses. h[ttps://github.com/molecularsets/moses](https://github.com/molecularsets/moses) (accessed
551    May 20, 2020).
552    (43) Johansson, S.; Prykhodko, O. Latent-Gan. h[ttps://github.com/Dierme/latent-gan](https://github.com/Dierme/latent-gan) (accessed Jan 5,
553    2021).
554    (44) Kotsias, P.; Bjerrum, E. J. Deepdrugcoder. h[ttps://github.com/pcko1/Deep-Drug-Coder](https://github.com/pcko1/Deep-Drug-Coder) (accessed Jan
555    5, 2021).
556    (45) Mercado, R.; Rastemo, T.; Lindelöf, E. Graphinvent. [https://github.com/MolecularAI/GraphINVENT/](https://github.com/MolecularAI/GraphINVENT/)
557    (accessed Oct 20, 2020).
558    (46) Rocío, M.; Tobias, R.; Edvard, L.; Günter, K.; Ola, E.; Hongming, C.; Esben Jannik, B., Graph Networks
559    for Molecular Design. *ChemRxiv* **2020**, *Preprint*.
560    (47) Oliphant, T. E., Python for Scientific Computing. *Computing in Science & Engineering* **2007**, *9*, 10-20.
561    (48) Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.;
562    Antiga, L. Pytorch: An Imperative Style, High-Performance Deep Learning Library. In Advances in neural
563    information processing systems, 2019; 2019; pp 8026-8037.
564    (49) Panaretos, V. M.; Zemel, Y., Statistical Aspects of Wasserstein Distances. *Annual review of statistics*
565    *and its application* **2019**, *6*, 405-431.
566    (50) Virtanen, P.; Gommers, R.; Oliphant, T. E.; Haberland, M.; Reddy, T.; Cournapeau, D.; Burovski, E.;
567    Peterson, P.; Weckesser, W.; Bright, J., Scipy 1.0: Fundamental Algorithms for Scientific Computing in
568    Python. *Nature methods* **2020**, *17*, 261-272.

569

570

571    Graphical TOC Entry



572