

1 **How good are current pocket based 3D generative models? :**
2 **The benchmark set and evaluation on protein pocket based**
3 **3D molecular generative models**

4

5 Haoyang Liu^{a,b,#}, Yifei Qin^{c#}, Zhangming Niu^{e,f,g,#}, Mingyuan Xu^b, Jiaqiang Wu^c,
6 Xianglu Xiao^{f,g,h}, Jinping Lei^{i,*}, Ting Ran^{b*}, Hongming Chen^{b,d,j*}

7

8 ^a State Key Laboratory of Medicinal Chemical Biology and College of Life Sciences,
9 Nankai University, 94 Weijin Road, Tianjin 300071, China

10 ^b Division of drug and vaccine research, Guangzhou National Laboratory, Guangzhou
11 315000, Guangdong, China

12 ^c School of pharmacy and food engineering, Wuyi University, Jiangmen 529020,
13 Guangdong, China

14 ^d School of Basic Medical Sciences, Guangzhou Laboratory, Guangzhou Medical
15 University, Guangzhou 511436, China

16 ^e National Heart and Lung Institute, Imperial College London, London SW7 2AZ, UK

17 ^f MindRank AI, Hangzhou, Zhejiang, China

18 ^g AI research center, MindRank Technologies Limited, London, UK

19 ^h Bioengineering Department and Imperial-X, Imperial College London, London W12
20 7SL, UK

21 ⁱ School of Pharmaceutical Science, Sun Yat-sen University, Guangzhou 510006,
22 China

23

24 ^j Lead contact

25

26 [#] Contributed equally

27

28 ^{*}Correspondence:

29 Hongming Chen (chen_hongming@gzlab.ac.cn); Ting Ran (ran_ting@gzlab.ac.cn);
30 [Jinping Lei \(leijp@mail.sysu.edu.cn\)](mailto:leijp@mail.sysu.edu.cn)

31

32 **Abstract**

33 The development of three-dimensional (3D) molecular generative model based on protein pockets
34 has recently attracted a lot of attentions. This type of model aims to achieve the simultaneous
35 generation of molecular graph and 3D binding conformation under the constraint of protein
36 binding. Various pocket based generative models have been proposed, however, currently there is

37 a lack of systematic and objective evaluation metrics for these models. To address this issue, a
38 comprehensive benchmark dataset, named as POKMOL-3D, is proposed to evaluate protein
39 pocket based 3D molecular generative models. It includes 32 protein targets together with their
40 known active compounds as a test set to evaluate the versatility of generation models to mimick
41 the real-world scenario. Additionally, a series of 2D and 3D evaluation metrics was integrated to
42 assess the quality of generated molecular structures and their binding conformations. It is expected
43 that this work can enhance our comprehension of the effectiveness and weakness of current 3D
44 generative models, and stimulate the discussion on challenges and useful guidance for developing
45 next wave of molecular generative models.

46

47 **Introduction**

48 Application of deep generative model in drug design has gained widespread attention. Over
49 the past few years, a large number of molecular generative models based on 1D/2D structures
50 have been reported. These models mainly generate molecules by learning the structural features
51 embodied in either 1D strings, such as the Simplified Molecular Input Line Entry System
52 (SMILES) strings¹ and SELFIE strings², or 2D molecular graphs^{3,4}. Despite substantial progress
53 being made in improving the validity of generated molecules and the efficiency of exploring the
54 drug-like chemical space, most of these models overlook the rich information contained in the 3D
55 conformation of molecules. Indeed, the binding affinity of a drug molecule to the target protein is
56 predominantly dependent on the degree of geometrical and electrostatic complementarity between
57 their 3D conformations. Therefore, in recent two years, molecular generative models based on 3D
58 conformations has become a hot research area.

59 Currently, 3D molecular generative models can be divided into three categories. One class
60 aims to generate 3D conformations of a given 2D molecule graph⁵⁻¹⁰. The second class is to
61 generate simultaneously 3D conformations and 2D graph of a molecule, which doesn't consider its
62 binding partner, i.e. the protein pocket¹¹⁻¹⁵. The third class strives to generate 3D conformations
63 and 2D graph under the constrain of protein binding pocket, which has attracted most interest in
64 the latest two years and is the ultimate goal of the so called structure based de novo molecular
65 design. In current study, we solely focus on the third type of 3D generative model.

66 The LiGAN model pioneered the field by introducing a conditional variational autoencoder
67 of 3D atomic density grids¹⁶, in which a protein pocket is encoded by a conditional encoder
68 network. The output atomic density grids are transformed to a 3D molecular structure using a
69 rule-based atom fitting and bond inference algorithm. The main drawback of this model lies in its
70 inability to maintain the equivariance on rotation and translation. In contrast, SBDD¹⁷ and
71 GraphBP¹⁸ leverage a 3D equivariant graph neural network (GNN) to solve this problem. They
72 both sequentially place atoms to a given 3D binding site, utilizing protein pocket and ligand atoms
73 generated in previous steps as the contextual information. For GraphBP model, an anchor atom
74 must be selected to determine position of the next generated atom, while SBDD generates the next
75 atom on a given arbitrary position in the pocket. However, they only generate atom types and
76 positions and utilize third-party software such as RDKit to construct bond types. Pocket2Mol¹⁹
77 introduces a geometric vector perceptron (GVP)-based equivariant GNN^{20, 21} to encode the 3D
78 geometric information of protein pocket and existing fragments of the ligand in the pocket.
79 Compared to GraphBP, it particularly involves a predictor network to infer the bond type. ResGen,
80 which utilizes similar GVP-based architecture, further encodes protein and ligand at residue and

81 atom levels respectively to better capture high-level binding interactions. Different from learning
82 the joint distribution of atom type and bond type in Pocket2Mol, ResGen decomposes the
83 distribution as a product of multiple conditional distributions for anchor atom, atom position, atom
84 type and bond type. PocketFlow²² adds a layer of geometric bottleneck perceptron (GBP) to the
85 GVP network to improve model speed and enhance information integration. It is also
86 characterized by its AtomFlow and BondFlow modules for predicting atom type and bond type,
87 respectively. Especially, chemical knowledge such as bond valence is explicitly integrated to guide
88 the bond inference. SurfGen²³ represents binding pocket as protein surface and utilizes a special
89 framework Geodesic-GNN to learn the distribution of the topological information on the surface.
90 All these GNN-based models can be categorized as autoregressive model, varying in the way of
91 encoding protein pocket and the decoding or sampling of atoms in the generative process.

92 Recently, diffusion model is an emerging deep learning technology utilizing an iterative
93 denoising process to map noise to data and have been used for 3D molecule generation.
94 DiffSBDD is the first 3D conditional graph diffusion model²⁴, in which protein pocket nodes
95 transformed from atomic point clouds are used as conditional constraints and remain unchanged
96 throughout the reverse diffusion process. TargetDiff is conceptually similar to DiffSBDD but
97 employs a different diffusion formalism for the categorical atom types. Both DiffSBDD and
98 TargetDiff map protein and ligand nodes into a joint embedding space for noise prediction, while
99 in DiffBP²⁵ these two types of node are separately embedded. In addition, DiffBP introduces a
100 new loss term to regulate the intersection between protein and ligand nodes in space. Besides,
101 language models have also been reported for 3D structure generation. Feng *et al.*²⁶ developed
102 Lingo3DMol that combines transformer-based language model architecture and deep geometric

103 learning technology for 3D molecular generation. A prior model was pre-trained to generate 3D
104 molecular structures given a fragment-based SMILES string, and then fine-tuned based on
105 protein-ligand complex data. The protein pocket and ligand embeddings are used as the input for
106 encoder and decoder respectively.

107 Despite various pocket-based 3D molecular generative models reported, there is still lack of
108 unified and comprehensive benchmark metrics to objectively evaluate the quality of generated
109 molecules. Early-developed models, such as GraphBP and SBDD, primarily relied on common
110 2D/3D molecular evaluation metrics such as molecular validity, molecular docking score,
111 druglikeness (QED)²⁷, synthesizability score (SAscore)²⁸, and structural diversity to assess the
112 quality of generated 3D molecules. Pocket2Mol additionally performed analysis on ring size of
113 molecules as part of quality measurement. Although TargetDiff, DiffSBDD and DiffBP were
114 published later than Pocket2Mol, they still adopted 2D molecular evaluation metrics. ResGen and
115 SurfGen introduced additional 2D metrics, e.g. the mean similarity between generated molecules
116 and known active molecules, to quantify their efficiency of generating active compounds.
117 Lingo3DMol analyzed the proportion of targets in which nearest neighbors of known active
118 compounds can be generated. In terms of measuring the quality of 3D conformation, the
119 Jensen-Shanon²⁹ divergences of bond length, bond angle, and dihedral angle of generated
120 molecules, and docking scores of redocked compounds in binding pocket are often used as the
121 metrics. While Lingo3DMol used “min-in-place” GlideScore³⁰ to evaluate the poses after
122 minimization of the generated conformations within the pocket. In addition, ResGen and SurfGen
123 employ extra 3D evaluation metrics such as *in situ* docking score, similarity of protein-ligand
124 interaction fingerprints between the generated conformation and known actives, and 3D similarity

125 index based on the overlay between the generated conformations and ground-truth conformations.

126 On the other hand, most of the pocket-based 3D molecular generative models were trained
127 and tested on the Crossdock2020 dataset³¹, which is constructed by molecular docking of active
128 ligands on PDB database to its corresponding targets. ResGen, SurfGen, and PocketFlow
129 additionally utilized the protein pockets outside the Crossdock dataset to assess the model's
130 potential for real-world application. In ResGen, two external pockets were selected for evaluation,
131 while in SurfGen the evaluation was extended to 20 therapeutic targets. In PocketFlow, the authors
132 synthesized two generated molecules for wet-lab validation, which experimentally validates the
133 effectiveness of the 3D generative models in hit finding scenario. Particularly, the resolved crystal
134 structures showed that the generated 3D conformations are highly similar to their active binding
135 conformations. So far, most of the pocket-based 3D generative models employ Pocket2Mol as the
136 baseline model for comparison, and the number of model included in their evaluation is relatively
137 small and incomplete. For future model development, it is probably necessary to conduct a
138 performance comparison among a larger model set under the same criteria. PoseCheck³² is a
139 small-scale benchmark study for this task by comparing five models including LiGAN,
140 Pocket2Mol, 3D-SBDD, Pocket2Mol, TargetDiff and DiffSBDD. PoseCheck focused on 3D
141 conformation evaluation using the CrossDock dataset as test set, and employed four 3D based
142 evaluation metrics, namely steric clashes based on van der Waals distance, protein-ligand
143 interaction fingerprints, strain energy of the generated conformations, and conformation similarity
144 between generated and docked poses. However, PoseCheck ignored 2D metrics that can also
145 imply the general quality of molecular structures. DrugPose is another small-scale benchmark
146 study focusing on 3D molecular generative models, but non-specific for the protein pocket based

147 methods³³. The binding similarity between pre- and post-docked poses of generated molecules,
148 drug-likeness and synthesizability were also analyzed. Zheng *et al.* recently conducted a
149 cross-algorithm benchmark study³² which compared a few protein pocket based 3D molecular
150 generative models with 1D SMILES/SELFIES and 2D molecular graph based generative methods.
151 Despite 16 models were evaluated in their study, only several commonly used metrics, such as
152 docking score, QED, molecular validity, were employed for the evaluation.

153 This study provides a comprehensive and systematic evaluation on nine 3D molecular
154 generative models in 32 protein pockets, and the compiled benchmark dataset is called
155 POKMOL-3D. In terms of evaluation metrics, both 2D and 3D metrics were considered and
156 classified according to their characteristics. Given that the essence of pocket-based 3D molecular
157 generation model is to generate molecules being able to bind specified targets, and the generated
158 conformations should be close to their active conformations, conventional 2D and 3D evaluation
159 metrics were expanded to include new parameters characterizing sampling speed, actives recovery
160 and conformation quality etc. Furthermore, the widely used SMILES based generative model
161 REINVENT^{34, 35} was included as the baseline model for comparison on the 2D based metrics,
162 providing an interesting perspective on how good current 3D based models comparing with
163 classical SMILES based model. In summary, this work could provide a systematic and
164 comprehensive benchmark set for evaluating 3D generative model.

165

166 **Method**

167 **Model selection**

168 Nine representative models were selected from recently published protein pocket-based 3D

169 molecular generative models spanning 2021 to 2024, which includes four distinct categories:
170 graph model, diffusion model, language model, and flow model. To demonstrate the efficiency of
171 3D generative model, the SMILES based REINVENT (version 4.0) was utilized as the base line
172 model.

Table 1. List of selected 3D molecular generative models

Model	Generative process	Model architecture	Training Set	Year
SBDD	Autoregressive	Graph Model	CrossDock2020	2021
GraphBP	Autoregressive	Graph model	CrossDock2020	2022
Pocket2Mol	Autoregressive	Graph model	CrossDock2020	2022
DiffBP	One-shot	Diffusion Model	CrossDock2020	2022
SurfGen	Autoregressive	Graph model	CrossDock2020	2023
TargeDiff	One-shot	Diffusion Model	CrossDock2020	2023
ResGen	Autoregressive	Graph model	CrossDock2020	2023
Lingo3DMol	Autoregressive	Language Model	PDBbind ³⁶	2024
PocketFlow	Autoregressive	Flow model	CrossDock2020	2024

173

174 **POKMOL-3D dataset**

175 In order to assess the versatility of selected models, 32 protein targets belonging to diverse
176 protein families were selected, in which five classes of target are included: kinases, non-kinase
177 enzymes, GPCRs, nuclear receptors, and protein-protein interaction targets. Given that our goal is
178 to evaluate model performance on generating molecules conditioned on the 3D information of
179 protein pocket, the targets that possess published protein-ligand complex structures were chosen.

180 For each target, one crystal structure from the RCSB PDB database³⁷, whose resolution is less than
181 3Å, was retrieved and only the subunit containing ligand was kept for analysis when the structure
182 comprises multiple subunits. All protein structures were optimized using the Protein Preparation
183 Wizard (PrepWizard) module in Maestro³⁸.

184 Moreover, active compounds of these proteins were extracted from the ChEMBL database³⁹
185 and served as reference set for evaluation metrics. Active molecules were considered eligible if
186 they have a molecular weight less than 500 Da, and the target IC₅₀/EC₅₀ values are less than 10
187 nM or Ki/Kd values less than or equal to 100 nM.

188 **Molecular generation**

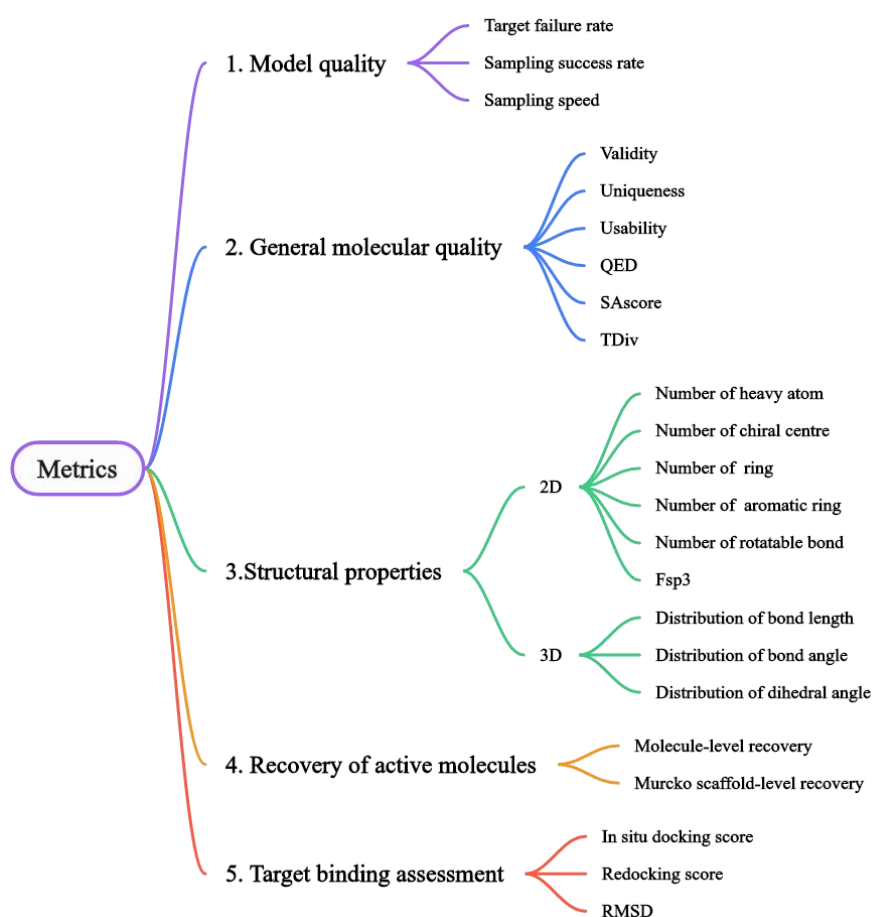
189 In this study, the latest version of the nine models were downloaded from GitHub. During the
190 sampling process, we adhered to the default configurations for all models, with the sole exception
191 of the sampling scale which was specifically calibrated to yield 2000 molecules per run. Each
192 model was tasked with generating a minimum of 2000 molecules per target. In scenarios where
193 sampling 2000 molecules in a single run was unfeasible for certain targets, maximal three
194 sampling runs were done to expand the generation set. For molecular generation employing the
195 REINVENT model, 1000 epochs of RL, steered by the molecular docking score (GlideScore),
196 were conducted to mimic the scenario of molecular generation within protein pocket. All
197 molecules generated during the RL process were subsequently utilized for further analysis.

198 **Evaluation metrics**

199 In current study, five categories of metrics were proposed to evaluate the performance of 3D
200 molecular generative models conditioned on protein pocket. As shown in Figure 1, it encompass
201 model quality, general molecular quality, structural properties, recovery of active molecules, and

202 target binding related scores. Specifically, the model quality class evaluates the sampling speed
203 and the target versatility of the model. The following two classes primarily focus on the quality of
204 molecular structure, providing insights into the overall molecular properties, 2D topological and
205 3D related properties of the generated molecules. The last two classes of metrics evaluate the
206 effectiveness of generating target binding compounds, at some extent reflecting the probability of
207 being able to bind the target protein. Through this pool of metrics, we hope to provide a thorough
208 benchmarking on current state-of-the-art 3D pocket-based molecular generative models, more
209 importantly providing guidance for developing new algorithms in future.

210



211

212 Figure 1. Five types of evaluation metrics for compared generative models.

213 **Model quality**

214 To evaluate the model quality, three metrics, *i.e.* target failure rate, sampling success rate and
215 sampling speed, were proposed. As shown in equation 1, target failure rate refers to the proportion
216 of targets for which no molecules can be generated by the model within three sampling runs.

$$\text{target failure rate} = \frac{\text{targets without molecules generated}}{\text{all targets}} \quad (1)$$

217 As shown in equation 2, sampling success rate refers to the proportion of targets for which the
218 model can generate more than 2000 molecules within three runs.

$$\text{sampling success rate} = \frac{\text{targets with more than 2000 molecules generated}}{\text{Total number of targets}} \quad (2)$$

219 In addition, the sampling speed was defined as the average time (in seconds) required for
220 generating one molecule. Here, the sampling time was counted for generating 100 molecules for
221 each target under the same computational resource. For calculating the sampling speed, the
222 employed computation resource was a linux workstation of 16-core 24GB RAM Intel Xeon
223 Platinum 8358 2.60GHz CPU and a NVIDIA GeForce 3090 GPU.

224

225 ***General molecular quality***

226 The general molecular quality set includes molecular validity, uniqueness, usability,
227 drug-likeness, synthetic score and target based diversity. These are properties reflecting overall
228 generation set. The calculation of properties was carried out using RDKit package⁴⁰. Molecular
229 validity, as defined in equation 3, refers to the proportion of valid molecules within the generated
230 set. Molecules that can successfully go through the standardization process are considered valid.

$$\text{validity} = \frac{\text{valid molecules}}{\text{generated molecules}} \quad (3)$$

231 As shown in equation 4, uniqueness refers to the proportion of unique molecules obtained
232 after removing duplicates among the valid molecules.

$$\text{uniqueness} = \frac{\text{unique molecules}}{\text{valid molecules}} \quad (4)$$

233 As shown in equation 5, usability refers to the proportion of molecules containing common
234 elements C, N, O, P, S, F, Cl, Br, I, and H. Molecules containing other elements are considered
235 unusable, for example those containing metal elements.

$$\text{usability} = \frac{\text{usable molecules}}{\text{unique molecules}} \quad (5)$$

236 Furthermore, the drug-likeness score (QED) and synthetic accessibility score (SAscore) were
237 calculated for each generated molecule using RDKit. These scores assess the drug-likeness and
238 synthetic feasibility of the generation set, respectively.

239 Target based molecular diversity (TDiv) is computed utilizing equation 6 to represent the
240 mean value of target specific diversity of generation sets:

$$\text{TDiv} = \frac{1}{N_{\text{target}}} \sum_{t=1}^{N_{\text{target}}} \left(1 - \frac{\sum_{i=1}^{N_t} \sum_{j=2}^{N_t} T_{\text{sim}}(i, j)}{N_t^2} \right), \quad i < j \quad (6)$$

241 where i and j refer to the indexes of two molecules in the generated molecule set for target t . The
242 Tanimoto similarity (T_{sim}) of the Morgan fingerprints⁴¹ is calculated based on all the pairs of
243 molecules for the same target. This similarity is normalized on the total number of molecules in
244 the generation set of the target. Target specific diversity score is derived from this normalized
245 similarity. The final TDiv score is defined as the average diversity score across all targets. A
246 higher TDiv value indicates greater diversity.

247

248 ***Structural properties***

249 The structural property group includes a set of 2D topological descriptors comprising the
250 number of heavy atom, chiral atom, ring, aromatic ring and rotatable bond, and the fraction of sp³
251 hybridized carbon atom (Fsp³), and their distribution was also compared. Additionally, a set of 3D

252 based geometrical properties containing the Jensen-Shanon divergency (JSD)²⁹ of bond length,
253 bond angle and dihedral angle was calculated. Detailed information for selected types of bond,
254 bond angle and dihedral angle can be found in Supporting Material Figure S1. The geometry
255 comparison was made between the conformations generated by the model and the low-energy
256 conformations optimized using the LigPrep module of Schrodinger package (version 2020). As
257 proposed in previous works²², JSD measures the distance of two probability distributions and is
258 defined as in Equations 7-8.

$$M = \frac{(P + Q)}{2} \quad (7)$$

$$JSD(P, Q) = \frac{1}{2}(D_{kl}(P||M) + D_{kl}(Q||M)) \quad (8)$$

259 where **P** denotes the probability distribution of a 3D property of the conformations generated by
260 the model, whereas **Q** corresponds to the conformations after energy optimization. **M** represents
261 the average distribution of P and Q. The Kullback-Leibler (KL) divergence⁴², denoted as D_{kl} , is
262 calculated separately to quantify the difference of either **P** or **Q** from M. The JSD value was then
263 obtained by averaging the KL divergences. A JSD value of 0 indicates that the distributions P and
264 Q are identical and a value of 1 represents completely dissimilar distributions.

265

266 *Recovery of active molecules*

267 Recovery of actives refers the ratio of generated compounds which are similar to the actives
268 in the reference set for a specific target. Tanimoto similarity between generation set and actives in
269 the reference set of the target protein is calculated. Here, Morgan fingerprint based on two bond
270 distance was used to calculate Tanimoto similarity. For an active molecule of target t (A_t), it is
271 recovered if its similarity of any compound in generation set is larger than 0.6. Then, the recovery

272 rate of active molecules for target t (R_t) is calculated as Formula 9:

$$R_t = \frac{\text{number of recovered active molecules}}{\text{total number of active molecules}} \quad (9)$$

273 This metric, at certain extent, can be regarded as the probability of reproducing active compound
274 by the generative model. The ratios at molecular structure and molecular Murcko scaffold⁴³ level
275 were examined respectively.

276 ***Target binding assessment***

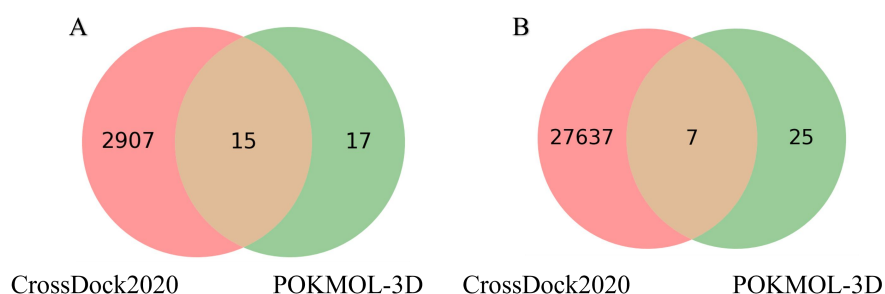
277 To evaluate how good the generated set can bind in its target pocket, two scoring strategies
278 were used here. One is the *in situ* scoring strategy, which scores the generated conformations at
279 the binding site without going through further pose optimization. The other one is the redocking
280 strategy, which scores after redocking of generated molecules into the binding set via external
281 docking software. The Glide docking module of Schrodinger software was used for *in situ* scoring
282 and redocking, and the GlideScore was used as the score value⁴⁴. In redocking, the generated
283 conformations were first gone through the LigPrep protocol of Schrodinger software for
284 preparation and then docked, and only one docking pose was saved for each molecule. Besides,
285 the Root-Mean-Square-Deviation (RMSD) value between the generated conformations and the
286 redocked conformations were also calculated without performing any conformation superposition.
287 The RMSD value quantifies the fitness between the protein pocket and the generated ligand, as
288 any Van der Waals clash or unmatched electrostatic interaction between ligand and protein would
289 penalize the ligand conformation, and in this case the docking pose was used as surrogate of
290 ground truth. In summary, the *in situ* score, redocking score and RMSD value between the
291 generated and redocked conformations were included in the target binding assessment set of
292 metrics.

293 Results and Discussion

294 Benchmark dataset composition

295 Protein pocket based 3D molecular generative model involves the encoding of 3D geometric
296 information of protein pocket, which are extracted from experimentally determined crystal
297 structures of various proteins. To investigate the generalizability of the investigated models to
298 unknown targets, our benchmark dataset encompasses 32 protein pockets, among which 17 targets
299 and 25 PDB IDs are not included in the CrossDock2020 dataset that usually used for training the
300 pocket based 3D generative models (Figure 2, detailed structure IDs can be seen in Table S1). In
301 addition, these targets belong to various druggable protein families, such as kinase, G
302 protein-coupled receptor (GPCR), nuclear receptor etc., and have reported ligands as marketed or
303 clinical drug.

304



305

306 **Figure 2.** Target overlap between the POKMOL-3D and CrossDock2020 datasets in terms of (A) protein target
307 name and (B) crystal structure.

308 Model quality

309

Table 2. Metrics for model quality evaluation

Model	Target failure rate	Sampling success rate	Sampling speed ^a
REINVENT	0	1	2.002

GraphBP	0	1	0.428
Pocket2Mol	0.094	0.656	7.131
PocketFlow	0	1	1.462
Lingo3DMol	0.063	0.313	2.381
DiffBP	0	1	16.338
TargetDiff	0	0.969	78.167
ResGen	0.063	0.125	29.613
SurfGen	0	0.094	28.150
SBDD	0	0.406	12.348

310 Note: a) estimated in second per compound

311 In this study, three structural unrelated metrics, i.e. target failure rate, sampling success rate
312 and sampling speed, were proposed to evaluate model quality. As shown in Table 2, most models
313 are able to generate compounds for all pockets so that their target failure rate is 0, while
314 Pocket2Mol, Lingo3DMol and ResGen fail to generate molecules for a few target proteins. In
315 detail, Pocket2Mol failed on Beta2AR, FXR and LXRB, ResGen failed on ERK2, NAMPT
316 proteins, and Lingo3DMol failed on CDK9 and DPP4 proteins. These results suggested that these
317 three models are not generalized good enough to deal with all targets. Additionally, the sampling
318 success rate, defined as the fraction of targets that a model can generate over 2,000 molecules at
319 most three runs, was employed as an additional indicator of generalizability to assess the models'
320 capacity to generate sufficient molecules given a specified sampling size. The results indicated
321 that models GraphBP, PocketFlow, DiffBP and the SMILES based baseline model REINVENT are
322 able to sample over 2,000 molecules for all targets within three sampling runs, and TargetDiff also
323 exhibits a high sampling success rate. Thus, the diffusion based and flow based models are able to
324 generate sufficient molecules from a model-type perspective. However, the remaining models, i.e.
325 Pocket2Mol, Lingo3DMol, ResGen, SurfGen, and SBDD, showed much lower success rate.
326 Notably, SBDD exhibited a significantly lower sampling success rate than GraphBP that share the
327 similar GNN architecture. This discrepancy might be attributed to the distinct approaches used by

328 these models to predict new atoms in the autoregressive generation process. A similar
329 phenomenon was observed when comparing ResGen with Pocket2Mol. SurfGen, which represents
330 protein pocket as protein surface, exhibited the lowest rate in sampling enough compounds in the
331 pockets. Although detailed reason is not unclear, one probable reason may be that the flaws
332 existed in generated 3D conformations make them failed in passing the internal structural validity
333 check.

334 Furthermore, comparison of the sampling speed was conducted. The results showed that
335 GraphBP exhibits fastest sampling speed, in which a molecule can be generated within one second.
336 In contrast, SBDD was much slower than GraphBP although they share similar generative
337 methodology. Pocket2Mol, PocketFlow, and Lingo3DMol exhibited relatively rapid sampling rate,
338 in which a compound can be sampled in less than 10 seconds. TargetDiff showed the slowest
339 speed, in which a compound is generated in more than one minute. Interestingly, DiffBP exhibited
340 much faster sampling speed than diffusion based TargetDiff and graph based models ResGen and
341 SurfGen. Furthermore, all protein pocket-based 3D molecular generative models, except GraphBP
342 and PocketFlow, exhibited slower sampling speed than REINVENT.

343 **General molecular quality**

344

Table 3. Comparison of general molecular quality

Model	Validity ↑	Uniqueness ↑	Usability ↑	QED ↑	SAscore ↓	Molecule Tdiv ↑	Scaffold Tdiv ↑
Actives	1	1	0.996	0.553	3.059	0.882	0.864
REINVENT	1	0.999	1	0.603	2.763	0.944	0.933
GraphBP	0.997	0.998	0.893	0.498	5.241	0.955	0.954

Pocket2Mol	0.906	0.903	1	0.420	4.074	0.926	0.925
PocketFlow	1	0.907	1	0.471	3.084	0.938	0.917
Lingo3DMol	1	0.742	1	0.484	3.17	0.913	0.871
DiffBP	0.992	0.998	1	0.508	3.553	0.942	0.935
TargetDiff	0.998	0.58	1	0.356	5.303	0.944	0.95
ResGen	1	0.992	1	0.345	4.169	0.927	0.922
SurfGen	1	0.987	1	0.369	4.408	0.927	0.923
SBDD	0.641	0.996	1	0.357	5.937	0.915	0.914

345

346 To assess the quality of 2D structures of generated molecules, six metrics were utilized:

347 molecular validity, uniqueness, usability, drug-likeness, synthesis accessibility, and diversity (as

348 shown in Table 3, the distribution plots can be seen in supporting material). Molecular validity and

349 uniqueness are fundamental metrics for evaluating generative models. The results indicated that

350 the molecular validity of all protein pocket-based 3D models except SBDD were either equal or

351 close to 1.0. In terms of molecular uniqueness, all models exhibited much higher performance

352 than Lingo3DMol and TargetDiff. Especially, a new metric named molecular usability was

353 introduced to quantify the likelihood of these models generating uncommon elements in structure.

354 The results indicated that all models exhibit good performance on this metric, and only GraphBP

355 generates about 10% molecules with uncommon atoms such as silicon atom.

356 In the assessment of drug-likeness, the QED score was averaged among the molecules

357 generated by each model (Figure S2A). Notably, the QED scores for most 3D models were below

358 0.5, lower than the average value observed for known active molecules (QED score = 0.553). In

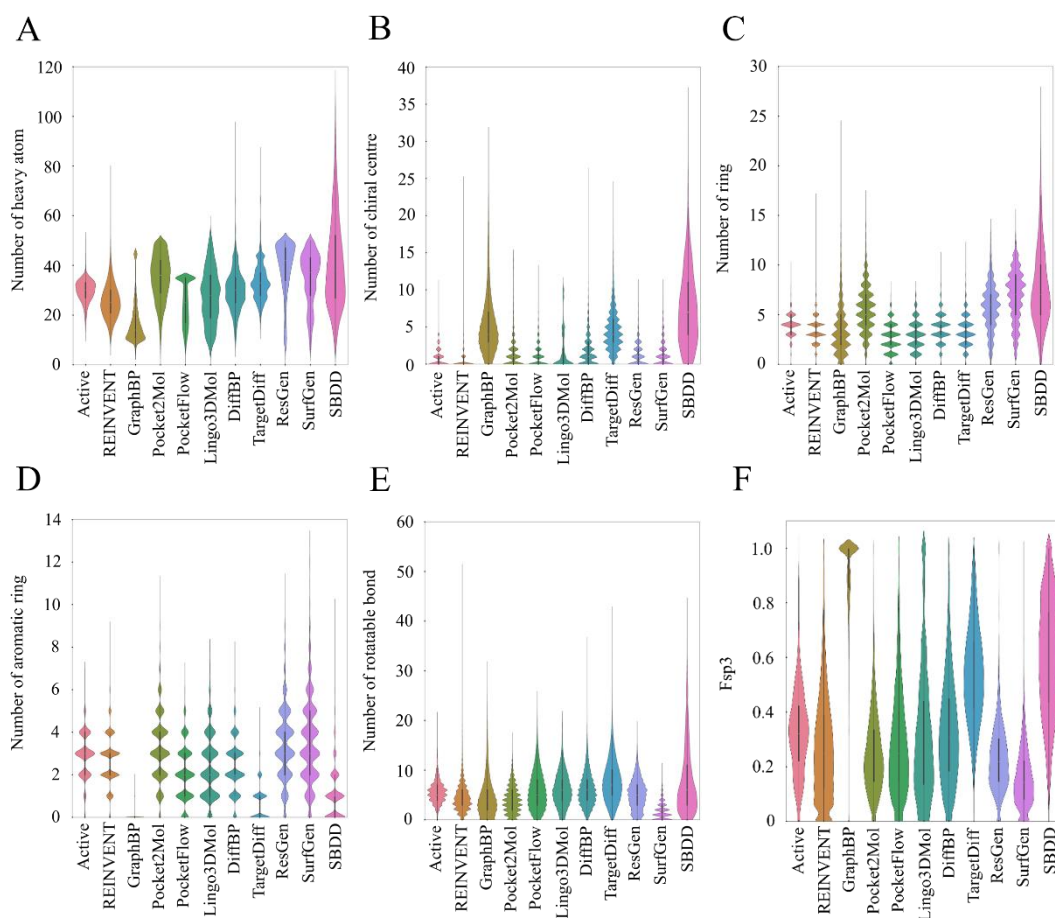
359 contrast, the QED score of REINVENT is superior to the average value of active molecules and
360 higher than all 3D models. Additionally, the SAScore metric was utilized to evaluate the synthetic
361 accessibility of generated molecules (Figure S2B). The SAScore of SBDD was significantly higher
362 than other models, suggesting that SBDD tended to generate molecules with poorer synthesis
363 accessibility. In contrast, PocketFlow and Lingo3DMol exhibited higher SAScore than other 3D
364 models but still worse than the baseline REINVENT.

365 In evaluating molecular diversity of generation set, a pairwise similarity calculation was
366 performed between molecules belonging to a specific target and the average molecular diversity
367 across all targets. The results listed in Table 3 revealed that both REINVENT and 3D generative
368 models exhibit high diversities. In summary, although 3D models exhibited similar performance
369 on validity, uniqueness, usability and diversity with baseline model REINVENT, REINVENT
370 model significantly showed better performance on QED and SAScore.

371

372 **Structural properties**

373 In addition to general molecular quality, it is imperative to consider fine-grained topology
374 related structural properties for evaluation. We analyzed the distribution of several crucial 2D
375 topological features, including the number of heavy atom, chiral centre, ring, aromatic ring, and
376 rotatable bond, and the Fsp3 (Figure 3A-F).



377

378 **Figure 3.** The violin plots for structural properties of 3D generative models including number of (A) heavy atom,
 379 (B) chiral centre, (C) ring, (D) aromatic ring, (E) rotatable bond, and (F) Fsp3.

380 The distribution of heavy atom number was shown in Figure 3A, REINVENT clearly
 381 exhibited most similar distribution as the active set, while most of 3D models tended to generate
 382 larger molecules than the active set, except GraphBP which generated significant portion of
 383 molecules with fewer than 20 heavy atoms. Particularly, SBDD generated molecules with a wide
 384 range from 20 to 60 heavy atoms. For Pocket2Mol, ResGen and SurfGen, the generated molecules
 385 had heavy atom between 40 and 60.

386 In terms of number of chiral centre, most of the 3D generative models tended to generate
 387 more chiral centre than the active set. Especially for GraphBP, SBDD and TargetDiff, the number
 388 of chiral centre in the generated compounds was obviously much larger than other models,

389 resulting in decrease of the synthetic accessibility. Whereas, REINVENT exhibited most similar
390 distribution to the active set and the generated molecules clearly had less chiral centre than the 3D
391 generative models. The ring count in a molecule serves as an indicator of its structural complexity,
392 given that most drug molecules possess at least one ring structure⁴⁵. The ring count distribution of
393 active set fell in the range of 3-5 rings. The REINVENT model showed most similar distribution
394 to the active set, whereas the 3D models mostly exhibited much broader distribution. Especially
395 for Pocket2Mol, ResGen, SurfGen and SBDD, a substantial proportion of molecules had more
396 than five rings, resulting the increases of structural complexity and decreases of drug-likeness. In
397 contrast, PocketFlow, Lingo3DMol, DiffBP and TargetDiff exhibited distributions more close to
398 the active set than other 3D models. Interestingly, GraphBP tended to generate compounds with
399 less rings than the active set, consistent with the observation on distribution of heavy atom count.

400 Given the prevalence of aromatic rings in drug molecules⁴⁶, the number of aromatic ring is
401 also an important metric. Figure 3D showed that the distribution of aromatic ring in REINVENT
402 was quite similar to that of the active set, in which most compounds have 2-4 aromatic rings,
403 while the distribution of 3D models was deviate from the active set. Interestingly, GraphBP,
404 TargetDiff and SBDD tended to generate compounds with less number of aromatic rings
405 comparing to the active set. Pocket2Mol, ResGen and SurfGen generated a substantial fraction of
406 compounds with more than four aromatic rings, while PocketFlow, Lingo3DMol, and DiffBP
407 generated compounds with primarily one to three aromatic rings.

408 The analysis on the number of rotatable bond (Figure 3E) revealed that REINVENT exhibits
409 most similar distribution to the active set, although it still had a minor fraction of molecules
410 exceeding ten rotatable bonds. Among the 3D generative models, Pocket2Mol and GraphBP were

411 similar to the active set, while others models generated larger fraction of compounds with more
412 than ten rotatable bonds, and SurfGen generated molecules with less than two rotatable bonds.

413 The fraction of sp^3 hybridized carbon atoms is a metric that partially reflects a molecule's
414 flatness, *i.e.* the larger fraction of aromatic ring in a molecule the smaller F_{sp3} value is, and it is
415 related to the success of drug in clinical trials⁴⁷. Typically, small-molecule oral drugs harbor
416 approximately 40% of their carbon atoms in the sp^3 hybridization state⁴⁸. Our findings (Figure 3F)
417 indicated that the proportion of sp^3 -hybridized carbon atoms in active molecules primarily fell
418 within the range of 20% to 40%. The molecules generated by REINVENT model exhibited a
419 distribution closely resembling that of the active set. The 3D models Pocket2Mol, PocketFlow,
420 Lingo3DMol, DiffBP, and ResGen exhibited comparable distributions to the active set (as shown
421 in Figure 3F). SurGen had a tendency of favoring molecules containing less than 20% sp^3
422 hybridized carbon atoms, indicating larger number of aromatic ring. Whereas, GraphBP,
423 TargetDiff and SBDD showed a large fraction of compounds with high F_{sp3} value, indicating
424 most of the carbon atoms in the structure are saturated carbons with few aromatic rings. The
425 analysis on the structural properties revealed that REINVENT has most close distribution to that
426 of the active set, while all the 3D generative models showed larger deviation to the active set,
427 highlighting the necessity of further improvement for current 3D generative model algorithms to
428 increase the compound quality.

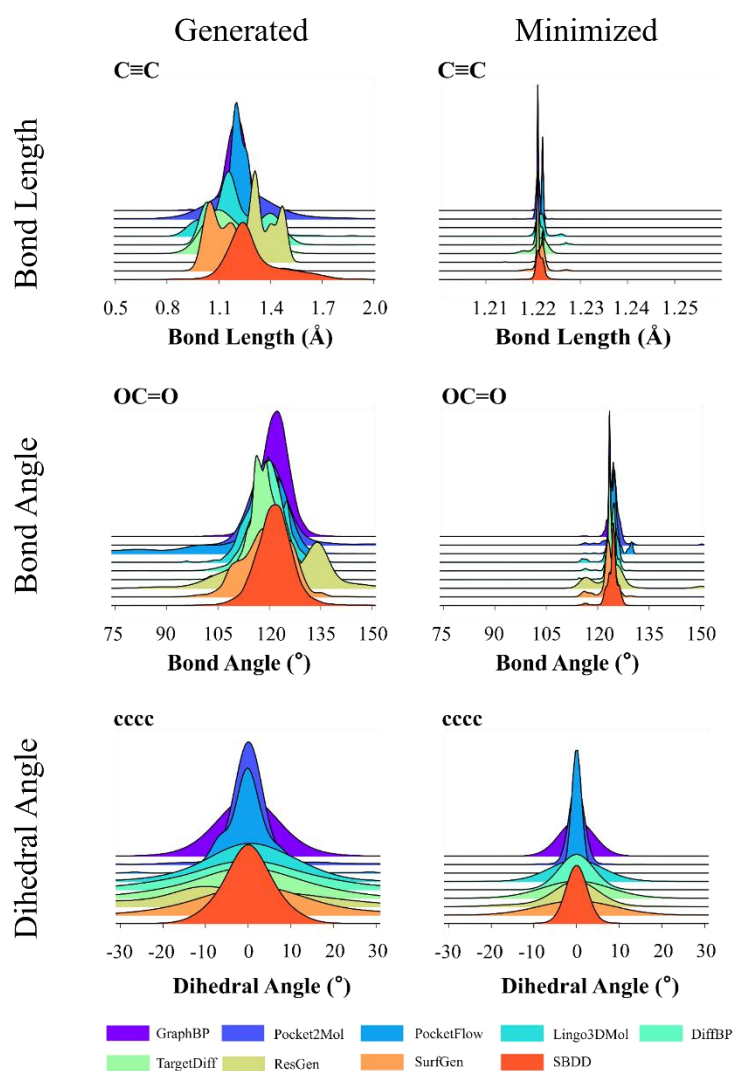
429 Besides the evaluation of 2D related metrics, the quality assessing of 3D conformation is also
430 important as these 3D based models generate 3D conformation and 2D graph simultaneously. We
431 investigated the differences of distributions of bond length, bond angle, and dihedral angle
432 between the generated 3D conformations and the OPLS3 force field ⁴⁹minimized conformations,

433 and the JSD index was utilized to quantify the deviation of these two distributions. For most 3D
434 models except GraphBP, the average JSD on bond length over all types of bonds is larger than 0.5
435 (Table 4), indicating significant difference in bond length distribution between the generated and
436 minimized conformations. Moreover, the JSD values for bond angles and dihedral angles surpass
437 0.1 for all models, indicating a great deviation to the force field minimized distributions, although
438 some 3D models are based on atomic point clouds leveraging third-party software such as
439 OpenBabel to construct the final structures. The distributions of $C\equiv C$ bond length, $OC=O$ bond
440 angle and $cccc$ dihedral angle were shown in Figure 4 (more detailed analysis in Figure S3-5), the
441 divergence could be attributed to the broader distribution on parameters in generated
442 conformations. Only a few bond angles and dihedral angles exhibited similar distribution
443 indicated by their small JSD metrics (Table S2-4). These results suggested that learning of bond
444 length, bond angle and dihedral angle in the generative models may need further improvement.
445

Table 4. JSD divergence of bond length, bond angle and dihedral angle

Model	Bond-Length↓	Bond Angle↓	Dihedral Angle↓
GraphBP	0.431	0.348	0.133
Pocket2Mol	0.586	0.404	0.145
PocketFlow	0.672	0.303	0.266
Lingo3Dmol	0.588	0.311	0.193
DiffBP	0.586	0.378	0.159
TargetDiff	0.610	0.367	0.180
ResGen	0.647	0.341	0.143

SurfGen	0.605	0.350	0.146
SBDD	0.653	0.348	0.238



446

447 **Figure 4.** Distribution of the bond length of C≡C, the bond angle of OC=O and the dihedral angle of cccc. The

448 left panel presents the distribution in generated conformations, while the right panel presents the minimized

449 distribution.

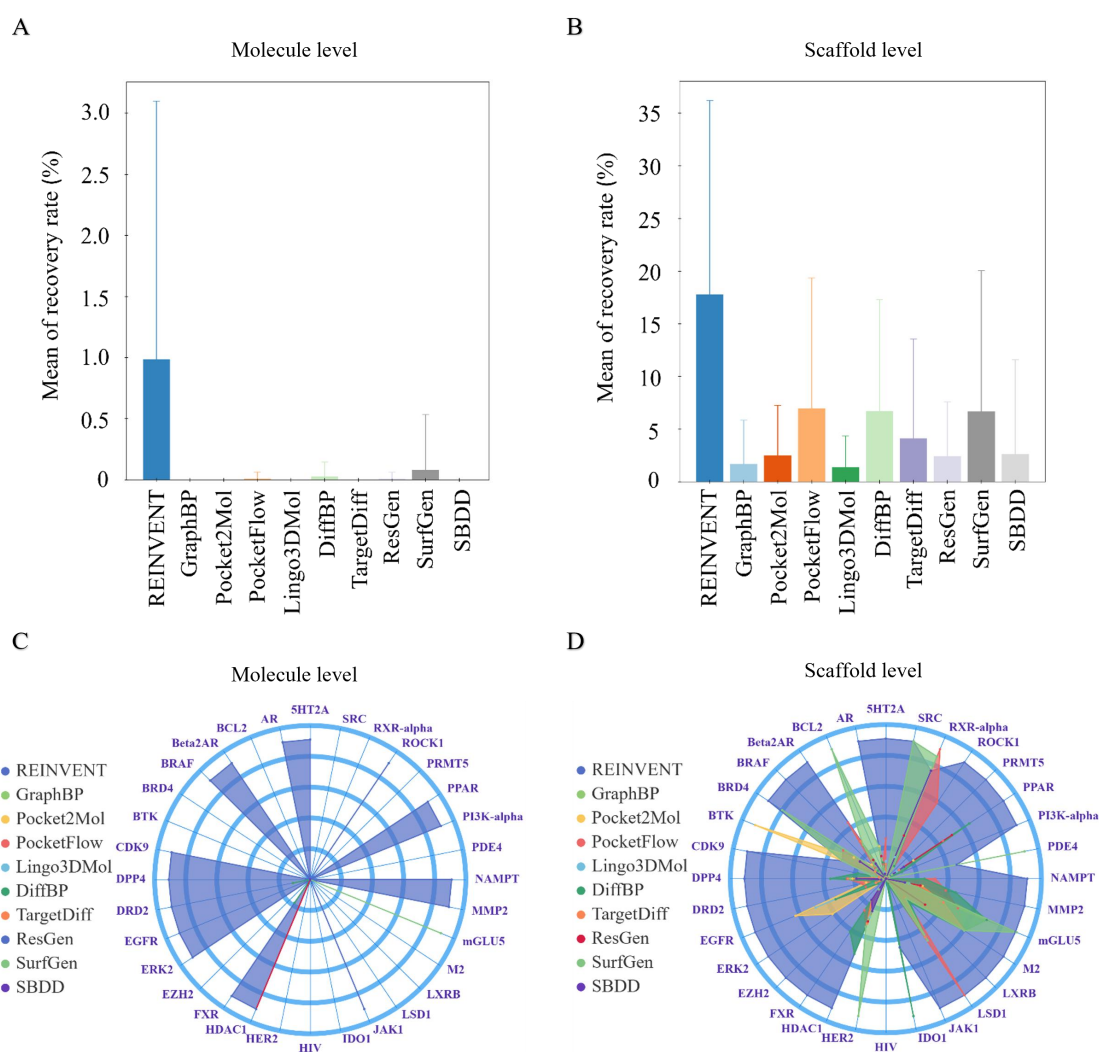
450 Recovery of active molecules

451 One direct way of evaluating 3D generative models conditioned on protein pockets is to

452 examine the model's capability of generating potentially bioactive molecules. For this purpose, we

453 employed the recovery rate of active molecules, *i.e.* the percentage of successfully recovered

454 active molecules for a given target. One active molecule was regarded as being successfully
 455 recovered if a similar compound in the generation set could be identified given the pair similarity
 456 was larger than the user defined cut-off. This is a stricter criterion compared to the similar metric
 457 utilized in Lingo3DMol²⁶, which is the percentage of targets with at least one generated molecule
 458 exhibiting similarity to the actives.



459

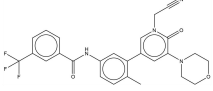
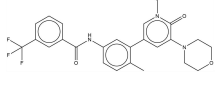
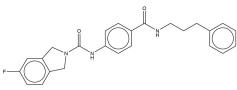
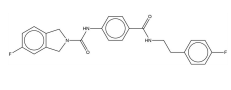
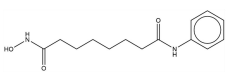
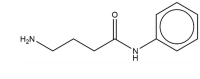
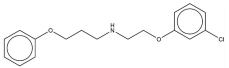
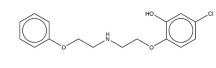
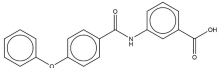
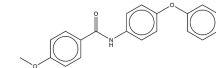
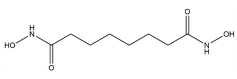
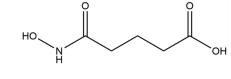
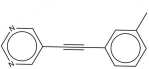
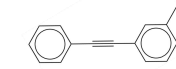
460 **Figure 5.** Statistics for recovery of active molecule. A-B) Histograms to display the average recovery rate

461 over all targets. C-D) Radar charts to display recovery rate among the target set for all models.

462 In current study, the Tanimoto similarity threshold was set to 0.6. The average recovery rates

463 at both molecule and scaffold level for these targets were shown in Figure 5A. It was obvious that

464 REINVENT showed the best recovery rate compared to all 3D molecular generative models at
 465 both levels, but there were still 15 targets in which REINVENT failed to recover any active
 466 molecule (Figure 5C). Specifically, GraphBP, Pocket2Mol, Lingo3DMol, TargetDiff and SBDD
 467 were unable to recover any active molecule for these 32 proteins (Table S5), while other 3D
 468 models can recovery a few actives on a few targets. The successfully recovered examples were
 469 shown in Figure 6. Compared to REINVENT, the 3D generative models tended to recover the
 470 active molecules with simple structure and relatively low similarity.

Model	Target	Active compound	Generate compound	Tanimoto similarity
REINVENT	BRAF			0.845
	NAMPT			0.782
PocketFlow	HDAC1			0.656
DiffBP	DRD2			0.610
	FXR			0.610
ResGen	HDAC1			0.650
SurfGen	mGLU5			0.621

471

472 **Figure 6.** Examples of successfully recovered active molecules for some models.

473

474 At Murcko scaffold level, a similar trend was observed among compared models (Figure 5B).

475 As anticipated, REINVENT exhibited superior performance (average recovery rate is around 17 %)

476 than the 3D models. PocketFlow, DiffBP, and SurfGen achieved relatively better performance than
477 other 3D models, but all 3D models exhibited average recovery rate lower than 10%. Furthermore,
478 significant variability was noted among the targets for each model. In terms of target coverage,
479 REINVENT clearly performed better on most targets (Figure 5D), while 3D models showed better
480 performance on a limited number of targets. For example, SurfGen performed better than
481 REINVENT on PDE4, HER2 and BCL2 proteins, Pocket2Mol on BTK and DiffBP on IDO1 etc.

482 Our results on active recovery rate demonstrated the limitation of current pocket based 3D
483 models as their performance was in general inferior than baseline REINVENT. Notably, even a
484 less stringent similarity threshold of 0.4 was applied, a similar trend was observed (Table S6),
485 indicating the weakness of these models in learning chemical structures conditioned on the protein
486 pocket.

487 **Target binding assessment**

488 The ultimate goal of generative model is to generate potential active compound to the target,
489 therefore we introduced several protein-ligand interaction based metrics to assess the target
490 binding capability of generated 3D conformations. Firstly, *in situ* docking score was chosen as a
491 surrogate to quantify the binding affinity of the 3D conformations generated from those 3D
492 models. Here, Glidescore was calculated for comparing binding affinity of conformations, and the
493 value of 0.0 kcal/mol was set as a criterion to judge if it is favorable for the ligand to bind in a
494 protein. A conformation is deemed to be a positive conformation (PC) if its Glidescore is less than
495 0.0 kcal/mol and the fraction of PC was calculated.

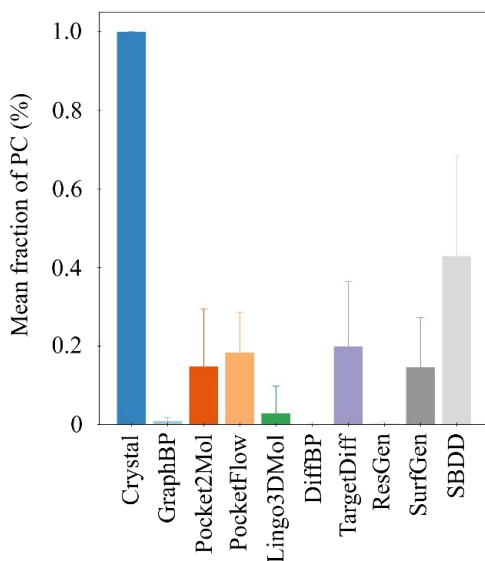


Figure 7. Histogram to display the mean fraction of PC over all targets.

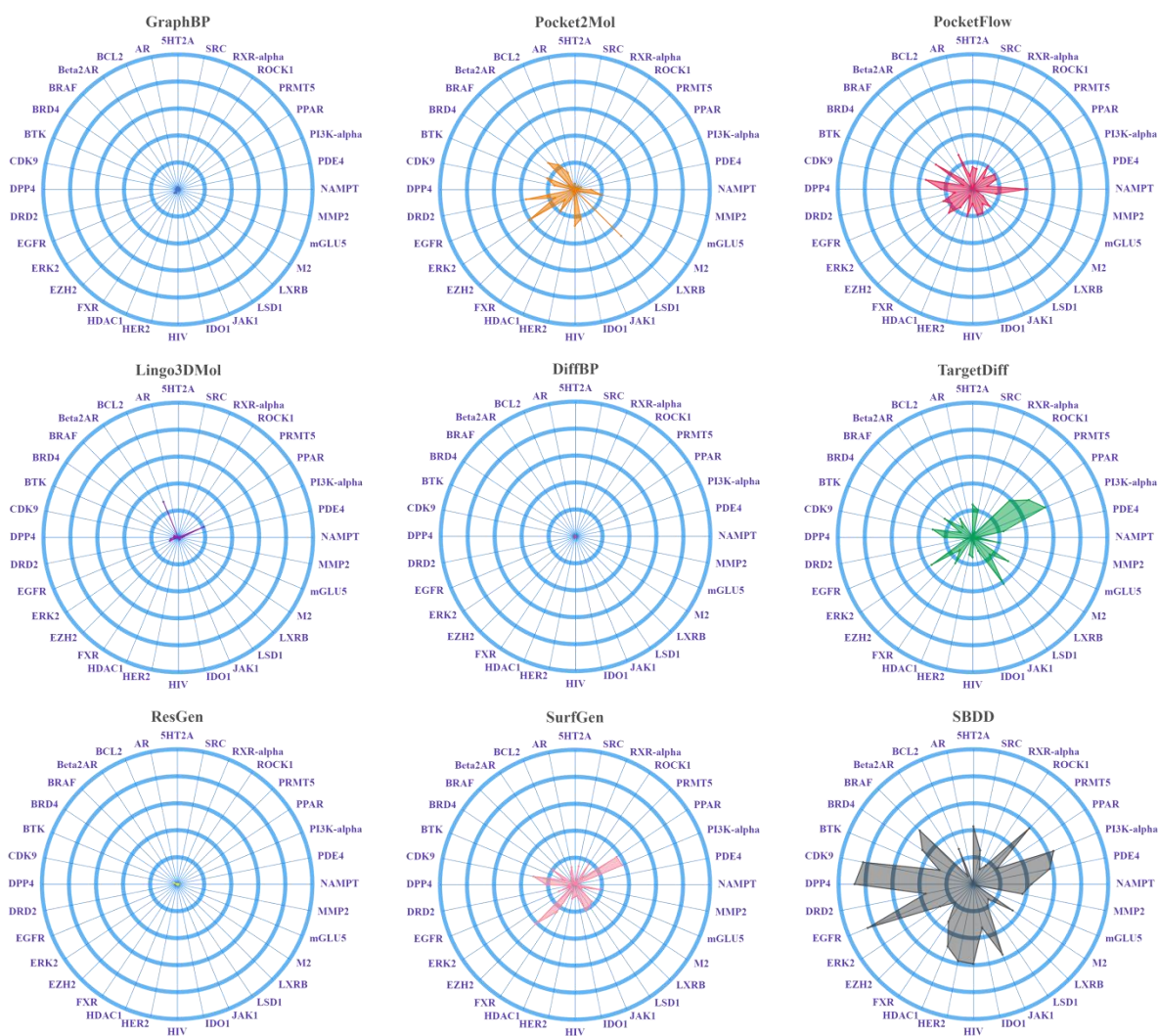


Figure 8. Radar charts to display the proportion of PC of each target across the models. The circles represent five

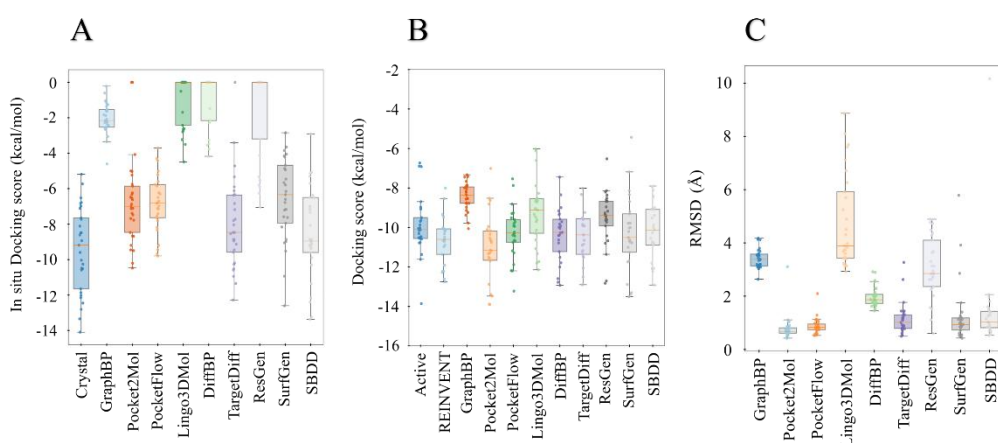
proportional levels, which are 100%, 80%, 60%, 40% and 20% from the outer layer to the centre of the chart.

501 The mean fraction of PC for each generation set could be seen in Figure 7 and the breakdown
502 of the fraction values across all 32 targets was shown in Figure 8. Among the 32 targets, GraphBP,
503 Lingo3DMol, DiffBP, and ResGen generated worst conformations and exhibited quite low
504 fraction values, while SBDD generated highest fraction of PC. A physically advantageous position
505 within the pocket was decided in SBDD model through a sampling process to estimate the
506 likelihood of atom occurrence and reduce the misplacement of ligand atoms in the pocket. SBDD
507 has on average only 40% conformations regarded as PC and only in very few targets the
508 proportion could surpass 80%. These results suggested that a lot of conformations have substantial
509 clash with protein atoms and current 3D models should strive to improve the learning of
510 protein-ligand interaction.

511 The distribution of *in situ* Glidescores for the top 10 conformations of 3D models, along with
512 that of crystal ligand conformations, was shown in Figure 9A. It can be seen from Figure 9A that
513 crystal ligands obviously exhibit best *in situ* docking scores, models SBDD, TargetDiff and
514 Pocket2Mol showed top three performances on *in situ* scores. Whereas, for Lingo3DMol, DiffBP
515 and ResGen with worst performances, the average *in situ* scores for some targets are even higher
516 than 0.0 kcal/mol. These results indicated that the learning of protein/ligand interaction is still far
517 from optimal for 3D generative models.

518 On the other hand, redocking analysis was also conducted to reproduce binding
519 conformations for the generated compounds via Glide docking in SP mode. The redocking
520 approach reflects the fitness between ligand and protein in 2D perspective as the binding
521 conformation and docking score is derived by external docking software. The distribution of
522 average redocking score for top 10 conformations was shown in Figure 9B, and scores of active

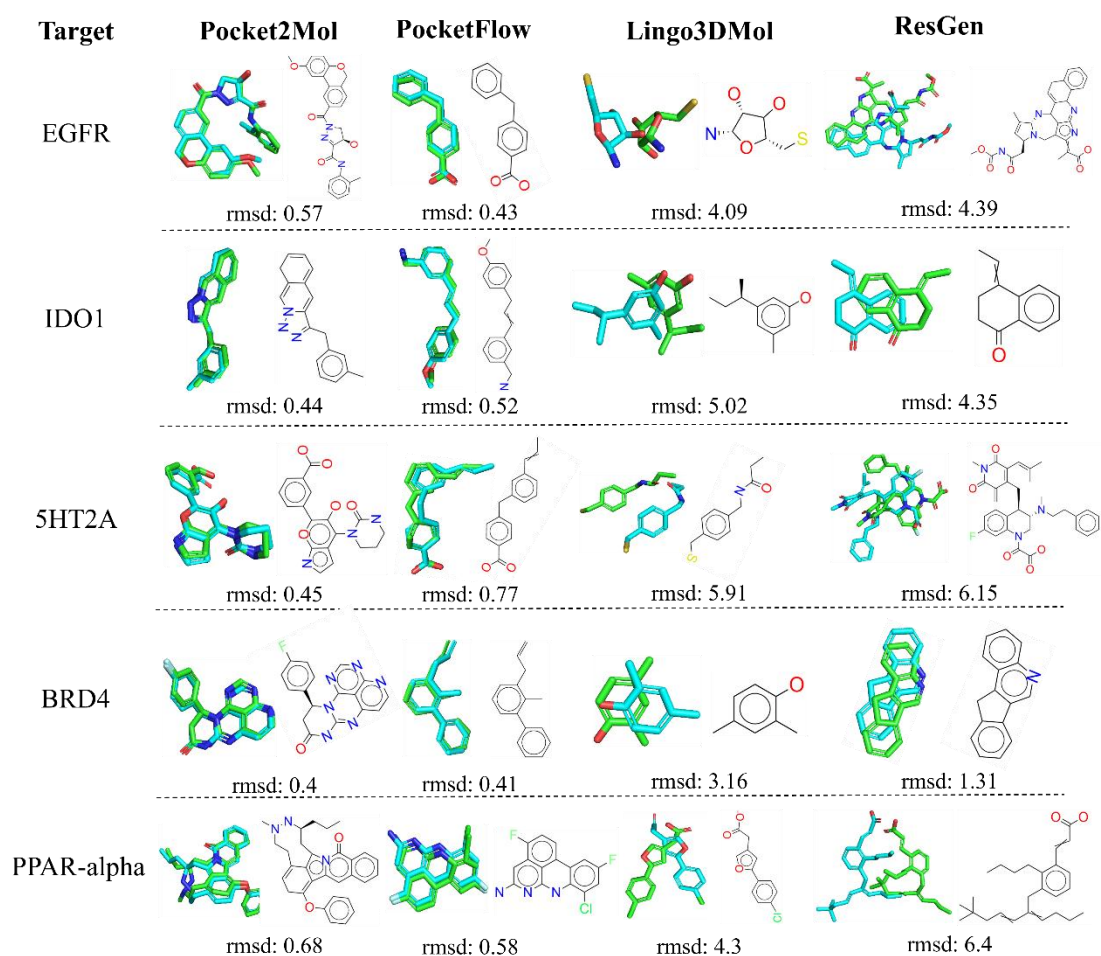
523 set and REINVENT were also included. The average redocking scores of all sets were between -8
524 and -12 kcal/mol, and differences of redocking scores among all models were small. In addition,
525 the average score of a few 3D models such as SBDD and Pocket2Mol were lower than the active
526 set, raising a question on whether the redocking score is really a relevant metric to evaluate
527 generative model.



528
529 **Figure 9.** Distribution of protein-ligand interaction based metrics across the targets calculated based on the top 10
530 molecules: A: in-situ score; B: redocking score; C: the RMSD between generated and redocked conformations.

531 It is obvious from Figure S6A-B that the redocking score is generally lower than the
532 corresponding in situ score, indicating that the generated conformation may be different from
533 active conformation. A comparative analysis was then conducted to measure the root mean square
534 deviation (RMSD) between the generated and redocked conformations (Figure S6C). As redocked
535 conformations were generated by force field based docking software, the RMSD value can
536 provide insights about geometrical difference between the generated 3D binding poses and poses
537 simulated by physics principles. The average RMSD of top 10 conformations was illustrated in
538 Figure 9C, which indicated that the RMSD values for the 32 targets fell in the range of [0, 10] Å.
539 ResGen, Lingo3DMol and GraphBP exhibited quite large RMSDs to the docked conformations,

540 while other models such as Pocket2Mol and PocketFlow showed average RMSD less than 2Å,
 541 suggesting the generated conformations of these models are close to their docked conformations.
 542 In Figure 10, several examples were presented to display overlapping between generated and
 543 docked conformations for the molecule with smallest RMSD for a specific target. As shown in
 544 Figure 10, Pocket2Mol and PocketFlow exhibited much better performance on all selected targets
 545 than Lingo3DMol and ResGen. The RMSD values of Pocket2Mol and PocketFlow were smaller
 546 than 1.0 Å, suggesting high similarity between the generated and docked poses. In contrast, the
 547 poses generated by Lingo3DMol and ResGen were quite dissimilar to the redocked poses in most
 548 cases.



549

550 **Figure 10.** Examples of overlapping between generated and redocked poses for Pocket2Mol, PocketFlow,

551 Lingo3Dmol, and ResGen. Five targets belonging to different families were selected for the comparison.

552 **Conclusion**

553 In current study, a novel benchmark dataset, named POKMOL-3D, was compiled
554 specifically for evaluating pocket based 3D molecular generative models, including a set of
555 comprehensive metrics for measuring model quality from various perspectives. Nine recently
556 published 3D generative models were selected for carrying out this benchmark study on 32 protein
557 pockets, along with the SMILES based REINVENT model as the baseline. Although some
558 promising results of pocket based 3D generative models have been reported, our benchmark study
559 revealed some weak points existed among the selected 3D models, such as slow sampling speed,
560 poor druggability and synthesizability of generated molecules, and failure to generate rational 3D
561 conformations for target binding. Overall, the performance of 3D generative models on large scale
562 of pockets is still far from satisfactory, and polishing of network architecture is needed to improve
563 the learning of ligand-protein interaction and generalizability of current 3D generative models to
564 enable their application on wide range of protein pockets. Through this study, we hope that the
565 proposed evaluation framework can be useful in facilitating the future advancement of pocket
566 based 3D molecular generative model.

567 **Conflicts of interest**

568 There are no conflicts to declare.

569

570 **Author contributions**

571 L. H. conceived the study, prepared the dataset and source code, and conducted the data analysis.

572 Q. Y. and Z. N. evaluated part of models and performed data validation. M. X., X. X and J. W. also

573 participate in discussion. C. H., R T., and L. J. are responsible for article revision and work
574 supervision.

575

576 **Acknowledgement**

577 We are grateful for support from the Big Data Center for Biomedical Research of Wuyi
578 University.

579

580 **Data availability**

581 The related data and source codes of POKMOL-3D are accessible at
582 <https://github.com/haoyang9688/POKMOL-3D>.

583

584 **Funding**

585 This study was supported by the National Key Research and Development Program of China
586 (2023YFF1204902), startup and R&D Program of Guangzhou National Laboratory
587 (YW-YWYM0205, GZNL2023A01008, GZNL2023A01005) and Overseas Experts Supporting
588 Programs under National Research Platform (WGZJ22-001).

589

590 **References**

- 591 1. D. Weininger, *Journal of Chemical Information and Computer Sciences*, 1988, **28**, 31-36.
- 592 2. A. Lo, R. Pollice, A. Nigam, A. D. White, M. Krenn and A. Aspuru-Guzik, *Digital Discovery*, 2023,
593 **2**, 897-908.
- 594 3. M. Skalic, D. Sabbadin, B. Sattarov, S. Sciabola and G. De Fabritiis, *Molecular pharmaceutics*,
595 **2019**, **16**, 4282-4291.
- 596 4. M. Xu, T. Ran and H. Chen, *Journal of Chemical Information and Modeling*, 2021, **61**,
597 3240-3254.
- 598 5. E. Mansimov, O. Mahmood, S. Kang and K. Cho, *Scientific reports*, 2019, **9**, 20381.

- 599 6. G. N. Simm and J. M. Hernández-Lobato, *arXiv preprint arXiv:1909.11459*, 2019.
- 600 7. T. Gogineni, Z. Xu, E. Punzalan, R. Jiang, J. Kammeraad, A. Tewari and P. Zimmerman,
601 *Advances in Neural Information Processing Systems*, 2020, **33**, 20142-20153.
- 602 8. O. Ganea, L. Pattanaik, C. Coley, R. Barzilay, K. Jensen, W. Green and T. Jaakkola, *Advances in*
603 *Neural Information Processing Systems*, 2021, **34**, 13757-13769.
- 604 9. S. Luo, C. Shi, M. Xu and J. Tang, *Advances in Neural Information Processing Systems*, 2021, **34**,
605 19784-19795.
- 606 10. C. Shi, S. Luo, M. Xu and J. Tang, International conference on machine learning, 2021.
- 607 11. L. Huang, H. Zhang, T. Xu and K.-C. Wong, Proceedings of the AAAI Conference on Artificial
608 Intelligence, 2023.
- 609 12. E. Hoogeboom, V. G. Satorras, C. Vignac and M. Welling, International conference on machine
610 learning, 2022.
- 611 13. A. Morehead and J. Cheng, *arXiv preprint arXiv:2302.04313*, 2023.
- 612 14. M. Xu, A. S. Powers, R. O. Dror, S. Ermon and J. Leskovec, International Conference on
613 Machine Learning, 2023.
- 614 15. C. Vignac, N. Osman, L. Toni and P. Frossard, Joint European Conference on Machine Learning
615 and Knowledge Discovery in Databases, 2023.
- 616 16. M. Ragoza, T. Masuda and D. R. Koes, *Chemical science*, 2022, **13**, 2701-2713.
- 617 17. S. Luo, J. Guan, J. Ma and J. Peng, *Advances in Neural Information Processing Systems*, 2021,
618 **34**, 6229-6239.
- 619 18. M. Liu, Y. Luo, K. Uchino, K. Maruhashi and S. Ji, *arXiv preprint arXiv:2204.09410*, 2022.
- 620 19. X. Peng, S. Luo, J. Guan, Q. Xie, J. Peng and J. Ma, International Conference on Machine
621 Learning, 2022.
- 622 20. B. Jing, S. Eismann, P. N. Soni and R. O. Dror, *arXiv preprint arXiv:2106.03843*, 2021.
- 623 21. C. Deng, O. Litany, Y. Duan, A. Poulencard, A. Tagliasacchi and L. J. Guibas, Proceedings of the
624 IEEE/CVF International Conference on Computer Vision, 2021.
- 625 22. Y. Jiang, G. Zhang, J. You, H. Zhang, R. Yao, H. Xie, L. Zhang, Z. Xia, M. Dai, Y. Wu, L. Li and S.
626 Yang, *Nature Machine Intelligence*, 2024, **6**, 326-337.
- 627 23. O. Zhang, T. Wang, G. Weng, D. Jiang, N. Wang, X. Wang, H. Zhao, J. Wu, E. Wang and G. Chen,
628 *Nature Computational Science*, 2023, **3**, 849-859.
- 629 24. A. Schneuing, Y. Du, C. Harris, A. Jamasb, I. Igashov, W. Du, T. Blundell, P. Lió, C. Gomes and M.
630 Welling, *arXiv preprint arXiv:2210.13695*, 2022.
- 631 25. H. Lin, Y. Huang, M. Liu, X. Li, S. Ji and S. Z. Li, *arXiv preprint arXiv:2211.11214*, 2022.
- 632 26. W. Feng, L. Wang, Z. Lin, Y. Zhu, H. Wang, J. Dong, R. Bai, H. Wang, J. Zhou and W. Peng,
633 *Nature Machine Intelligence*, 2024, **6**, 62-73.
- 634 27. G. R. Bickerton, G. V. Paolini, J. Besnard, S. Muresan and A. L. Hopkins, *Nature chemistry*,
635 2012, **4**, 90-98.
- 636 28. P. Ertl and A. Schuffenhauer, *Journal of cheminformatics*, 2009, **1**, 1-11.
- 637 29. J. Lin, *IEEE Transactions on Information theory*, 1991, **37**, 145-151.
- 638 30. R. A. Friesner, J. L. Banks, R. B. Murphy, T. A. Halgren, J. J. Klicic, D. T. Mainz, M. P. Repasky, E.
639 H. Knoll, M. Shelley and J. K. Perry, *Journal of medicinal chemistry*, 2004, **47**, 1739-1749.
- 640 31. P. G. Francoeur, T. Masuda, J. Sunseri, A. Jia, R. B. Iovanisci, I. Snyder and D. R. Koes, *Journal of*
641 *chemical information and modeling*, 2020, **60**, 4200-4215.
- 642 32. C. Harris, K. Didi, A. R. Jamasb, C. K. Joshi, S. V. Mathis, P. Lio and T. Blundell, *arXiv preprint*

643 *arXiv:2308.07413*, 2023.

644 33. Z. Jocys, J. Grundy and K. Farrahi, *Digital Discovery*, 2024.

645 34. M. Olivecrona, T. Blaschke, O. Engkvist and H. Chen, *Journal of cheminformatics*, 2017, **9**,
646 1-14.

647 35. H. H. Loeffler, J. He, A. Tibo, J. P. Janet, A. Voronov, L. H. Mervin and O. Engkvist, *Journal of*
648 *Cheminformatics*, 2024, **16**, 20.

649 36. R. Wang, X. Fang, Y. Lu and S. Wang, *Journal of medicinal chemistry*, 2004, **47**, 2977-2980.

650 37. H. M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T. N. Bhat, H. Weissig, I. N. Shindyalov and P.
651 E. Bourne, *Nucleic acids research*, 2000, **28**, 235-242.

652 38. S. Maestro, *Schrödinger, LLC, New York, NY*, 2020, **2020**, 682.

653 39. B. Zdrzil, E. Felix, F. Hunter, E. J. Manners, J. Blackshaw, S. Corbett, M. de Veij, H. Ioannidis, D.
654 M. Lopez and J. F. Mosquera, *Nucleic acids research*, 2024, **52**, D1180-D1192.

655 40. G. Landrum, *Release*, 2013, **1**, 4.

656 41. D. Rogers and M. Hahn, *Journal of chemical information and modeling*, 2010, **50**, 742-754.

657 42. S. Kullback, 1951.

658 43. G. W. Bemis and M. A. Murcko, *Journal of medicinal chemistry*, 1996, **39**, 2887-2893.

659 44. R. A. Friesner, R. B. Murphy, M. P. Repasky, L. L. Frye, J. R. Greenwood, T. A. Halgren, P. C.
660 Sanschagrin and D. T. Mainz, *Journal of medicinal chemistry*, 2006, **49**, 6177-6196.

661 45. R. D. Taylor, M. MacCoss and A. D. Lawson, *Journal of medicinal chemistry*, 2014, **57**,
662 5845-5859.

663 46. T. J. Ritchie and S. J. Macdonald, *Drug discovery today*, 2009, **14**, 1011-1020.

664 47. F. Lovering, J. Bikker and C. Humblet, *Journal of medicinal chemistry*, 2009, **52**, 6752-6756.

665 48. D. C. Kombo, K. Tallapragada, R. Jain, J. Chewning, A. A. Mazurov, J. D. Speake, T. A. Hauser
666 and S. Toler, *Journal of chemical information and modeling*, 2013, **53**, 327-342.

667 49. E. Harder, W. Damm, J. Maple, C. Wu, M. Reboul, J. Y. Xiang, L. Wang, D. Lupyan, M. K.
668 Dahlgren and J. L. Knight, *Journal of chemical theory and computation*, 2016, **12**, 281-296.

669