

---

# PoseCheck: Generative Models for 3D Structure-based Drug Design Produce Unrealistic Poses

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Deep generative models for structure-based drug design (SBDD), where molecule  
2 generation is conditioned on a 3D protein pocket, have received considerable inter-  
3 est in recent years. These methods offer the promise of higher-quality molecule  
4 generation by explicitly modelling the 3D interaction between a potential drug  
5 and a protein receptor. However, previous work has primarily focused on the  
6 quality of the generated molecules themselves, with limited evaluation of the 3D  
7 *poses* that these methods produce, with most work simply discarding the generated  
8 pose and only reporting a “corrected” pose after redocking with traditional meth-  
9 ods. Little is known about whether generated molecules satisfy known physical  
10 constraints for binding and the extent to which redocking alters the generated  
11 interactions. We introduce POSECHECK, an extensive analysis of multiple state-  
12 of-the-art methods and find that generated molecules have significantly more  
13 physical violations and fewer key interactions compared to baselines, calling into  
14 question the implicit assumption that providing rich 3D structure information im-  
15 proves molecule complementarity. We make recommendations for future research  
16 tackling identified failure modes and hope our benchmark will serve as a spring-  
17 board for future SBDD generative modelling work to have a real-world impact.  
18 Our evaluation suite is easy to use in future 3D SBDD work and is available at  
19 <https://anonymous.4open.science/r/posecheck-358E>.

## 20 1 Introduction

21 Structure-based drug design (SBDD) [1, 2, 3] is a cornerstone of drug discovery. It uses the 3D  
22 structures of target proteins as a guide to designing small molecule therapeutics. The intricate atomic  
23 interactions between proteins and their ligands shed light on the molecular motifs influencing binding  
24 affinity, selectivity, and drug-like properties. Employing computational methods such as molecular  
25 docking [4, 5], molecular dynamics simulations [6], and free energy calculations [7], SBDD aids in  
26 the identification and optimization of potential drug candidates.

27 Deep generative models for SBDD have recently attracted considerable attention in the ML commu-  
28 nity [8, 9]. These models learn from vast compound databases to generate novel chemical structures  
29 with drug-like properties [10]. By explicitly integrating protein structure information, these models  
30 aim to generate ligands with a higher likelihood of binding to the target protein. In particular, advance-  
31 ments in geometric deep learning [11, 12, 13] have led to a new suite of generative methods, enabling  
32 the design of 3D molecules directly within the confines of the target protein [14, 15, 16, 17, 18].  
33 These methods, which concurrently generate a molecular graph and 3D coordinates, provide the  
34 significant advantage of obviating the need for determining the 3D pose *post hoc* through traditionally  
35 slow molecular docking programs – at least in theory.

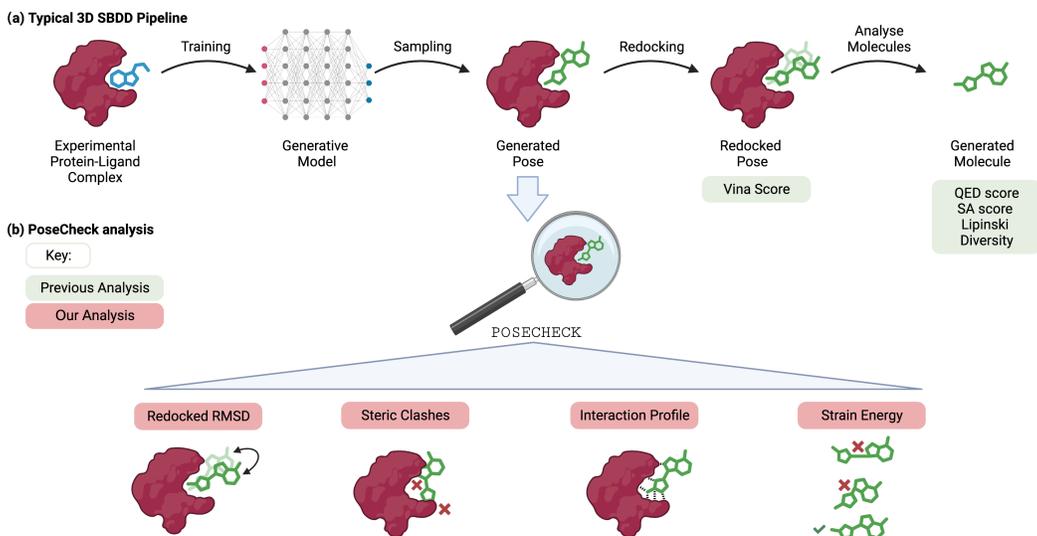


Figure 1: **Top:** Overview of a conventional pipeline of SBDD with 3D generative modelling. A generative model is usually trained using experimental or synthetic protein-ligand complexes, from which new molecules and poses can be sampled *de novo*. Typically, generated poses are discarded and redocked into the receptor, and primarily evaluated on 2D molecular graphs (e.g. QED). The effect of redocking on the final complex is often unknown, preventing understanding of the common failure modes of the trained model and therefore inhibiting progress. **Bottom:** the POSECHECK benchmarks for generated poses include pipeline-wide as well as component-wise metrics, enabling a targeted evaluation of each model component guiding further model development.

36 Assessing the quality of molecules generated by these methodologies is not straightforward, with  
 37 little work on experimental validation, especially for *de novo* design [19]. Typical evaluation metrics  
 38 (Figure 1a) focus primarily on the 2D graph of the generated molecules themselves, measuring  
 39 their physicochemical properties (e.g. QED [20]) and adherence to drug discovery heuristics (e.g.  
 40 Lipinski’s Rule of Five [21]). For effective SBDD, we argue that it’s equally essential to assess the  
 41 quality of the generated *binding poses* and their capacity to satisfy known biophysical prerequisites  
 42 for binding (Figure 1b). This perspective is essential if these methods are to serve as practical  
 43 alternatives to traditional virtual screening approaches in SBDD.

44 We hypothesise that multiple failure modes, undetected by currently applied metrics, are inherent  
 45 within these methods. The situation is further complicated by the common practice of disregarding  
 46 the initially generated pose and then redocking the molecule to attain a potentially enhanced pose  
 47 and only reporting these scores [16, 15, 18, 22]. This strategy tends to focus on presenting only  
 48 the outcomes of the redocked molecule, and it is not clear whether molecules shown in figures are  
 49 generated or redocked, making the accurate assessment of pose quality an increasingly intricate  
 50 challenge.

51 Our primary contributions are summarized as follows: We introduce POSECHECK, a set of new  
 52 biophysical benchmarks for SBDD models, expanding the traditional ‘pipeline-wide’ framework by  
 53 integrating ‘component-wise’ metrics (i.e. generated and redocked poses), leading to comprehensive  
 54 and precise model assessment. Utilizing this new framework, we evaluate a selection of high-  
 55 performing machine learning SBDD methods, revealing two key findings: (1) generated molecules  
 56 and poses often contain nonphysical features such as steric clashes, hydrogen placement issues,  
 57 and high strain energies, and (2) redocking masks many of these failure modes. Based on these  
 58 evaluations, we propose targeted recommendations to rectify the identified shortcomings. Our work  
 59 thus provides a roadmap for addressing critical issues in SBDD generative modelling, informing  
 60 future research efforts.

## 61 2 Background

62 **Deep Generative Models for 3D Structure-based Drug Design** Many works have recently tried  
63 to recast the SBDD problem as learning the 3D conditional probability of generating molecules given  
64 a receptor, allowing users to sample new molecules completely *de novo* inside a pocket. Common  
65 methods utilize Variational AutoEncoders (VAEs) [23], Generative Adversarial Networks (GANs)  
66 [24], Autoregressive (AR) models and recently Denoising Diffusion Probabilistic Models (DDPMs)  
67 [25]. LiGAN [14] uses a 3D convolutional neural network combined with a VAE model and GAN-  
68 style training. 3DSBDD [15] introduced an autoregressive (AR) model that iteratively samples from  
69 an atom probability field (parameterised by a Graph Neural Network) to construct a whole molecule,  
70 with an auxiliary network deciding when to terminate generation. Pocket2Mol [16] extended this  
71 work with a more efficient sampling algorithm and better encoder. DiffSBDD [17], DiffBP [26] and  
72 TargetDiff [17] are all conditional DDPMs conditioned on the 3D target structure. DecompDiff [27]  
73 is another diffusion model that decomposes the ligand into fragments for which it considers separate  
74 priors for the diffusion process. FLAG [22] chooses a fragment from a motif vocabulary based on  
75 the protein structure and composes it with other motifs into a final ligand in an iterative fashion.  
76 GraphBP [28] utilises an autoregressive flow model to formulate the ligand design as a sequential  
77 generation task.

78 **Related work** Guan et al. [17] perform limited analysis of small chemical sub-features, such as  
79 agreement to experimental atom-atom distances and the correctness of aromatic rings within the  
80 generated molecule. Baillif et al. [19] emphasize the necessity of 3D benchmarks for 3D generative  
81 models. However, both of these works study the molecules in isolation rather than the protein-ligand  
82 context. Both DecompDiff [27] and DiffBP [26] take steric clashes into account via their loss  
83 functions, but do not include steric clashes as a metric in their evaluation. TargetDiff [17] includes an  
84 analysis of Vina Scores but does not report any standard deviations on these. However, these standard  
85 deviations are critical in evaluating the performance of these models as we demonstrate in this paper.

86 The concurrent work PoseBusters [29] also focuses on benchmarking the biophysical plausibility of  
87 protein-ligand poses but focuses on evaluating *docking tools* instead of molecular generation models.  
88 They also find generalisation to new sequences to be poor.

## 89 3 Methods

90 In order to evaluate the quality of generated poses and their capacity to facilitate high-affinity  
91 protein-ligand interactions, we present a variety of computational methods and benchmarks in this  
92 section. These methodologies provide a thorough perspective on the poses produced and illuminate  
93 the ability of generative models to generate trustworthy and significant ligand conformations. Full  
94 implementation details are given in Appendix A.

95 **Interaction fingerprinting** Interaction fingerprinting is a computational method utilized in SBDD  
96 to represent and analyze the interactions between a ligand and its target protein. This approach  
97 encodes specific molecular interactions, such as hydrogen bonding and hydrophobic contacts, in a  
98 compact and easily comparable format – typically as a bit vector, known as a *interaction fingerprint*  
99 [30, 31]. Each element in the interaction fingerprint corresponds to a particular type of interaction  
100 between the ligand and a specific residue within the protein binding pocket. We compute interactions  
101 using the ProLIF library [30].

102 **Steric clashes** In the context of molecular interactions, the term *steric clash* is used when two  
103 neutral atoms come into closer proximity than the combined extent of their van der Waals radii [32].  
104 This event indicates an energetically unfavourable [33], and thus physically implausible, interaction.  
105 The presence of such a clash often points towards the current conformation of the ligand within the  
106 protein being less than optimal, suggesting possible inadequacies in the pose design or a fundamental  
107 incompatibility in the overall molecular topology. Hence, the total number of clashes serves as a vital  
108 performance metric in the realm of SBDD. We stipulate a clash to occur when the pairwise distance  
109 between a protein and ligand atom falls below the sum of their van der Waals radii, allowing a clash  
110 tolerance of 0.5 Å.

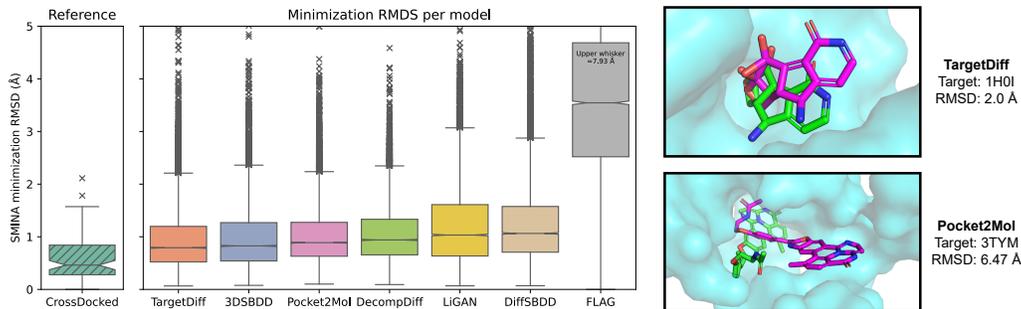


Figure 2: **Left:** RMSD between the generated and SMINA minimized poses for CrossDocked and all generative methods (note FLAG upper whisker value is not shown to preserve a meaningful scale). **Right:** Examples of large conformational rearrangements in the ligand upon redocking.

111 **Strain-energy** Strain energy refers to the internal energy stored within a ligand as a result of  
 112 conformational changes upon binding. When a ligand binds to a protein, both the ligand and the  
 113 protein may undergo conformational adjustments to accommodate each other, leading to changes  
 114 in their bond lengths, bond angles, and torsional angles. These changes can cause strain within the  
 115 molecules, which can affect the overall binding affinity and stability of the protein-ligand complex  
 116 [34]. Whilst there is always a trade-off between enthalpy and entropy, generally speaking, lower strain  
 117 energy results in more favourable binding interactions and potentially more effective therapeutics.  
 118 We calculate the strain energy as the difference between the internal energy of a relaxed pose and  
 119 the generated pose (without pocket). Both relaxation and energy evaluation are computed using the  
 120 Universal Force Field (UFF) [35] using RDKit.

121 **Docking** Our final assessment involves measuring the level of agreement between the docking  
 122 programs and the molecules produced by the learned distribution in the generative model. Although  
 123 this is the most coarse-grained approach we employ and docking programs come with their inherent  
 124 limitations, they nevertheless contain useful proxies and serve as valuable tools for comparison. In  
 125 this procedure, we redock the generated pose using SMINA [36]. Next, we compute the Root Mean  
 126 Squared Deviation (RMSD) between the generated pose and the docking-predicted one across all  
 127 generated molecules, thereby obtaining a distribution of RMSD values.

## 128 4 Results

### 129 4.1 Experimental Setup

130 In our study, we evaluate the quality of poses from seven recent methods: LiGAN [14], 3DSBDD  
 131 [15], Pocket2Mol [16], TargetDiff [17], DiffSBDD [18], DecompDiff [27] and FLAG [22]. All  
 132 models were trained on the CrossDocked2020 [37] dataset using the dataset splits computed in Peng  
 133 et al. [16], which used a train/test split of 30% sequence identity to give a test set of 100 target  
 134 protein-ligand complexes which we use for evaluation. For each model, we sampled 100 molecules  
 135 per target. We give a more detailed overview of the CrossDocked dataset and its limitations in  
 136 Appendix A.

### 137 4.2 Agreement with docking scoring functions

138 **Results** To discern whether the generated poses/binding modes produced by these models corre-  
 139 spond to overall low energy states with few physical violations, our preliminary analysis involves  
 140 determining the extent to which minimized poses preserve information from the initially generated  
 141 binding mode. Therefore, we proceed to compute the RMSD between the model-generated pose and  
 142 SMINA-minimized pose [36], with a lower RMSD value denoting a higher degree of agreement.<sup>1</sup>

143 The distributions of SMINA-minimization RMSDs of various methods are illustrated in Figure 2.  
 144 We first consider CrossDocked as a baseline, which has a mean minimization RMSD of 0.59 Å.

<sup>1</sup>To provide perspective, it's worth noting that a carbon-carbon bond generally measures 1.54 Å in length.

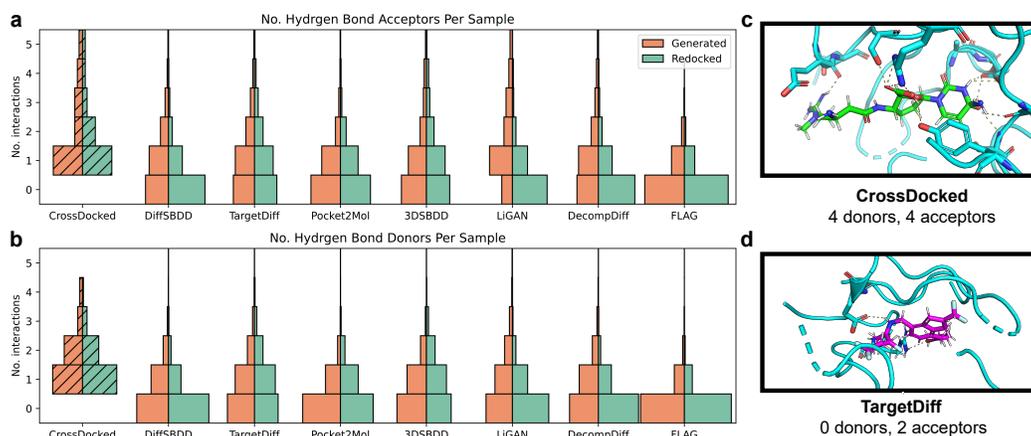


Figure 3: Interactions between protein and ligands as seen in generated poses (orange) and redocked poses (green). The frequency of (a) hydrogen bond acceptors and (b) hydrogen bond donors are considered. We find that generative models have significant trouble making hydrogen bond interactions compared to baseline (shaded boxes). Vertical histogram box sizes are normalised along the x-axis such that all have the same area. (c) Example from CrossDocked with large hydrogen bonding network. (d) Typical example from a generative models with low number of HBs.

145 Given that all the generative models were trained on these poses, we would expect to observe similar  
 146 performance. However, we find that all methods (except FLAG) have a mean score between 0.94 and  
 147 1.28 Å, suggesting that the generated binding poses are very far from low-energy states. We observe  
 148 little correlation between method types here except for the two similar AR models, 3DSBDD and  
 149 Pocket2Mol, which obtain mean RMSDs of 0.99 and 1.02 Å respectively. FLAG is the most egregious  
 150 example with on average 3.64 Å RMSD during minimization and a maximum value of 10.72 Å, an  
 151 extreme value for local minimisation.

152 **Discussion** These findings raise concerns for several reasons. They expose the minimal concordance  
 153 between the binding models learned by these methods and the established SMINA methodology  
 154 [5], despite it being the source of training data. More critically, they underline the lack of accurate  
 155 evaluations of generative models' capability to produce realistic binding poses; instead, these models  
 156 tend to generate drug-like molecules with vague binding modes, later rectified through docking.

157 We also calculated the RMSD between the generated and highest affinity redocked pose but were not  
 158 able to discern any reasonable signal-to-noise over the baseline dataset. We hypothesise that this may  
 159 be due to the fact that Francoeur et al. [37] provided up to 20 poses for every ligand, resulting in 22.5  
 160 million complexes, and the processing done in Peng et al. [16] is not clear on which poses they chose,  
 161 meaning these models may not have been trained on the lowest affinity poses.

### 162 4.3 Protein-ligand interaction analysis

163 **Evaluation** Below describe the classes of interaction that we evaluate. **Hydrogen bonds (HBs)** are  
 164 a type of interaction that occurs between a hydrogen atom that is bonded to a highly electronegative  
 165 atom, such as nitrogen, oxygen, or fluorine [38]. They are key to many protein-ligand interactions  
 166 [39] and require very specific geometries to be formed [40]. The directionality of HBs confers unique  
 167 identities upon the participating atoms: hydrogen atoms attached to electronegative elements are  
 168 deemed 'donors', whilst the atom accepting the HB is termed an 'acceptor'. **Van der Waals contacts**  
 169 (vdWs) are interactions that occur between atoms that are not bonded to each other. These forces can  
 170 be attractive or repulsive and are typically quite weak [41]. However, they can be significant when  
 171 many atoms are involved, as is typical in protein-ligand binding [42]. **Hydrophobic interactions**  
 172 are non-covalent interactions that occur between non-polar molecules or parts of molecules in a  
 173 water-based environment. They are driven by the tendency of water molecules to form hydrogen  
 174 bonds with each other, which leads to the exclusion of non-polar substances. This exclusion principle  
 175 prompts these non-polar regions to orient away from the aqueous environment and towards each other  
 176 [43], thereby facilitating the association between protein and ligand molecules [44].

177 **Results** Distributions of hydrogen bonding interactions are shown in Figure 3. We consider whether  
 178 our generative models can design molecules with adequate hydrogen bonding and find that no method  
 179 can match or exceed the baseline. In the reference set, CrossDocked, the modal number of HBs for  
 180 both acceptors and donors is 1, with means of 2.23 and 1.66 for acceptors and donors respectively.  
 181 Strikingly, we find that in all generated poses for all models (except LiGAN HB acceptors) the *most*  
 182 *common number of HB acceptors and donors is 0*, with means varying between 0.36-1.73 for HB  
 183 acceptors and 0.26-0.85 for HB donors. We find an average difference of 0.50 and 0.81 HBs between  
 184 the best-performing models and the baseline for acceptors and donors respectively. Results for Van  
 185 der Waals contacts and hydrophobic interactions are closer to the dataset baseline (see Appendix  
 186 Figure 6), possibly as these are easier to form.

187 **Discussion** Conventional wisdom would suggest that many minor imperfections in the generated  
 188 pose would be simply fixed by redocking the molecule (e.g. moving an oxygen atom slightly to  
 189 complete a hydrogen bond.) We find this is in fact rarely the case, with redocking sometimes being  
 190 significantly deleterious (see examples of LiGAN in Figure 3), suggesting that there are either  
 191 limitations in the docking function used or, more likely, the generated interaction was physically  
 192 implausible to begin with.

#### 193 4.4 Clash scores

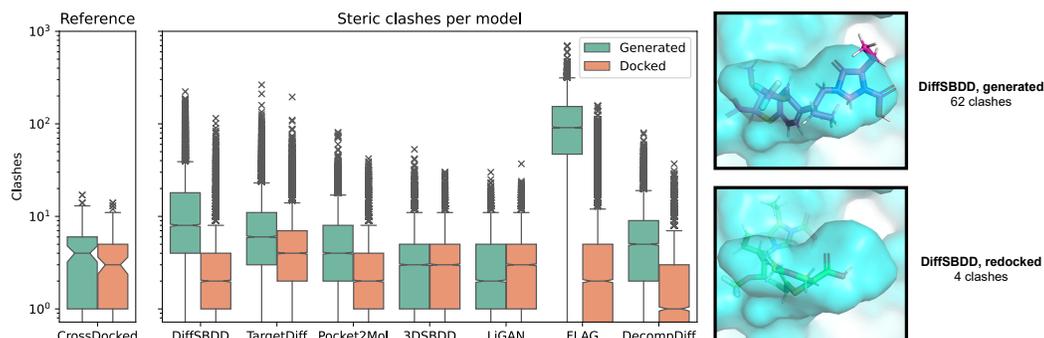


Figure 4: **Left:** number of steric clashes for the CrossDocked reference dataset as well as for the molecules generated by each model, both before and after redocking. **Right:** examples of a generated pose (magenta) and the same pose after redocking (green).

194 **Results** Figure 4 presents the results of the steric clash analysis. In summary, the latest methods,  
 195 particularly those employing diffusion models and fragment libraries, exhibit poor performance in  
 196 terms of steric clashes compared to the baseline, with a significant number of outliers. Although  
 197 redocking mitigates clashes to some extent, it does not always resolve the most severe cases.

198 The CrossDocked test set has a low number of clashes with few extreme examples, with a mean  
 199 of 4.59, upper quantile of 6 and maximum value of 17. In terms of generated poses, the older  
 200 methods perform best, with 3DSBDD and LiGAN having means of 3.79 and 3.40 clashes respectively.  
 201 Pocket2Mol, an extension of 3DSBDD, performs worse with a mean clash score of 5.62 and upper  
 202 quantile of 8 clashes. Finally, the diffusion-based approaches perform poorly with mean clash scores  
 203 of 15.33, 9.03 and 7.13 for DiffSBDD, TargetDiff and DecompDiff respectively. The tail end of their  
 204 distributions is also high, with the methods having upper quantiles of 18, 11 and 9 clashes respectively,  
 205 with TargetDiff having the worst case of 264 steric clashes. FLAG has the worst generated clash  
 206 scores by far, with mean and median clash scores of 110.96 and 91 respectively. Redocking the  
 207 molecules generally fixed many clashes and improved scores, especially for FLAG, where the mean  
 208 clash score improves from 110.96 to 5.55. The mean clash score for Pocket2Mol improves from 5.62  
 209 to 2.98, TargetDiff from 9.08 to 5.79 and DiffSBDD from 15.34 to 3.61.

210 **Discussion** Interestingly, DiffSBDD and TargetDiff, which are considered state-of-the-art based  
 211 on mean docking score evaluations [17, 18], exhibit subpar performance in their number of clashes.  
 212 They aim to learn atom position distributions without explicit constraints on final placements. While  
 213 DiffSBDD starts with a performance deficit, its enhanced clash mitigation during redocking elevates  
 214 its results to match the baseline, highlighting methodological distinctions between it and TargetDiff.

215 Notably, 3DSBDD and LiGAN show low clash scores, with the former positioning atoms within a  
216 predefined voxel grid [15] and the latter applying a clash loss [14]. DecompDiff also applies a steric  
217 clash loss (but does not directly measure clashes in the corresponding publication) [27] and performs  
218 best out of all the diffusion-based approaches. Generated molecules for FLAG were most egregious  
219 here; we speculate this is a result of first choosing a fragment from a fragment vocabulary using a  
220 softmax function and then forcing the placement of the fragment [22], regardless of whether it fits  
221 sterically.

222 Our findings affirm the assumption that redocking alleviates many minor clashes, akin to the force-  
223 field relaxation step in AlphaFold2 [45]. We initially speculated that molecules with clashes exceeding  
224 100 had been mistakenly generated inside the protein pocket. Yet, we often discovered fragments  
225 within highly constrained nooks, especially worsened with the addition of hydrogen atoms.

226 **Limitations** An important consideration to bear in mind is that proteins are not entirely rigid  
227 receptors. They can often experience limited conformational rearrangements to accommodate  
228 molecules of varying shapes and sizes [46]. Consequently, conducting generation and redocking in a  
229 rigid receptor environment may not yield accurate scores for potentially plausible molecules. Note all  
230 these results are with a *generous* clash tolerance of 0.5 Å (roughly half the vdW radii of a hydrogen  
231 atom), in order to be able to resolve differences between methods.

## 232 4.5 Strain energy

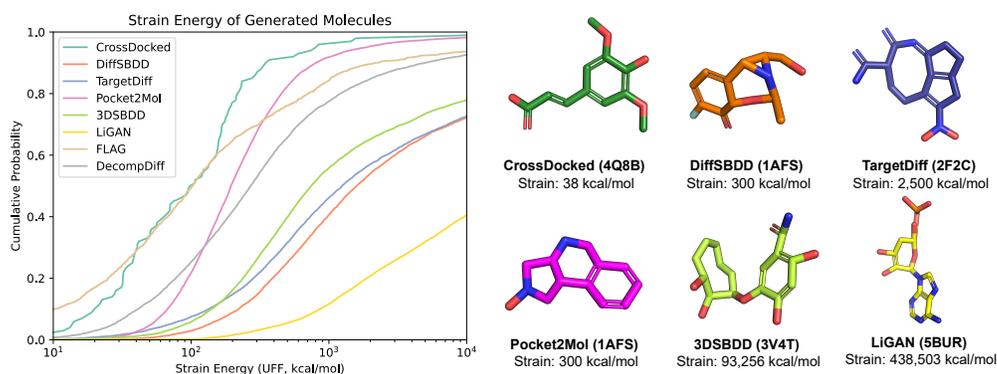


Figure 5: **Left:** CDF of strain energies. **Right:** Examples of molecules with high strain energy.

233 **Results** To conclude our study, we provide an analysis of the strain energy [34] of the generated  
234 poses. Force field relaxation before docking is a common post-processing step of many generative  
235 SBDD pipelines, masking some potential issues with the generated geometries less clear. This allows  
236 us to evaluate the generated molecules for undesirable properties like unrealistic bond distance or  
237 impossible geometries in rings.

238 Figure 5 displays the cumulative density function (CDF) of strain energy for the generated molecules,  
239 with the CrossDocked dataset serving as a baseline (Note: the x-axis is on a logarithmic scale). We  
240 focus on median values in our discussion since they are more representative in this context due to the  
241 presence of extreme outliers, with *mean* values ranging from approximately 10<sup>4</sup> to 10<sup>15</sup> kcal/mol.  
242 None of the generative methods yields molecules exhibiting strain energy close to that of the test set,  
243 which has a median strain energy of 102.5 kcal/mol.

244 **Discussion** Intriguingly, both of the diffusion-based methodologies (DiffSBDD and TargetDiff)  
245 perform similarly poorly, reporting median values of 1243.1 and 1241.7 kcal/mol, respectively. This  
246 could suggest issues with the currently used noised schedules [47] of these methods for ultra-precise  
247 atom position refinement (discussed in Section C). 3DSBDD performs to the same order of magnitude,  
248 with a median strain energy of 592.2 kcal/mol, suggesting that placing atoms into a discretized voxel  
249 space [15], while good for avoiding clashes, has a detrimental impact on the strain energy.

250 FLAG performs the best by far here with a median of 101.1 kcal/mol. We believe this due to most  
251 of the bond angles and distances already consisting of idealised geometries when the fragments are

252 initialized for incorporation into the molecule. Out of the other methods, Pocket2Mol performs  
253 the best in terms of strain energy, with a median of 194.9 kcal/mol. The method provides perhaps  
254 the finest-grained control over exact coordinates generated, by first choosing a focal atom and then  
255 generating a new atom coordinate directly using an equivariant neural network [13, 16], which may  
256 allow for more precise placement. LiGAN exhibits the highest strain energy, with a median value of  
257 18693.8 kcal/mol, indicating the poorest performance in this context.

258 **Limitations** The exceedingly high strain energy values observed in this scenario should be ap-  
259 proached with considerable prudence. For comparison, the combustion of TNT releases approximately  
260 815 kcal/mol. [48]. This data is not to be perceived as absolute, but rather illustrative of the extent to  
261 which our generated geometries deviate markedly from the standard distribution for the force field.  
262 This further underscores the existing issues. It is also conceivable that these poses might not even be  
263 initialized within more sophisticated, high-fidelity force fields [49].

## 264 5 Conclusion

265 In conclusion, this work presents a comprehensive exploration of structure-based drug design (SBDD)  
266 methodologies with deep generative models. We advocate for the need to consider *both* the quality  
267 of the generated molecules *and* the quality of the binding poses in these models, calling for an  
268 expanded evaluation of SBDD. The application of deep generative models in SBDD holds promise for  
269 developing innovative drug-like molecules. However, for SBDD approaches to realise that potential,  
270 we must establish a rigorous evaluation regimen of both the generated molecules and their interaction  
271 with the target – as proposed in this paper. Our research provides a solid evaluation regimen for future  
272 advancements in this field and we hope that it stimulates further development towards more efficient  
273 drug discovery processes.

## 274 References

- 275 [1] Tom L Blundell. Structure-based drug design. *Nature*, 384(6604 Suppl):23–26, 1996.
- 276 [2] Leonardo G Ferreira, Ricardo N Dos Santos, Glaucius Oliva, and Adriano D Andricopulo.  
277 Molecular docking and structure-based drug design strategies. *Molecules*, 20(7):13384–13421,  
278 2015.
- 279 [3] Amy C Anderson. The process of structure-based drug design. *Chemistry & biology*, 10(9):  
280 787–797, 2003.
- 281 [4] Oleg Trott and Arthur J Olson. Autodock vina: improving the speed and accuracy of docking  
282 with a new scoring function, efficient optimization, and multithreading. *Journal of computational*  
283 *chemistry*, 31(2):455–461, 2010.
- 284 [5] Amr Alhossary, Stephanus Daniel Handoko, Yuguang Mu, and Chee-Keong Kwoh. Fast,  
285 accurate, and reliable molecular docking with quickvina 2. *Bioinformatics*, 31(13):2214–2216,  
286 2015.
- 287 [6] John L Klepeis, Kresten Lindorff-Larsen, Ron O Dror, and David E Shaw. Long-timescale  
288 molecular dynamics simulations of protein structure and function. *Current opinion in structural*  
289 *biology*, 19(2):120–127, 2009.
- 290 [7] Christophe Chipot and Andrew Pohorille. *Free energy calculations*, volume 86. Springer, 2007.
- 291 [8] Yuanqi Du, Tianfan Fu, Jimeng Sun, and Shengchao Liu. Molgensurvey: A systematic survey  
292 in machine learning models for molecule design. *arXiv preprint arXiv:2203.14500*, 2022.
- 293 [9] Clemens Isert, Kenneth Atz, and Gisbert Schneider. Structure-based drug design with geometric  
294 deep learning. *Current Opinion in Structural Biology*, 79:102548, 2023.
- 295 [10] Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato,  
296 Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel,  
297 Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven  
298 continuous representation of molecules. *ACS central science*, 4(2):268–276, 2018.

- 299 [11] Michael M Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst.  
300 Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, 34  
301 (4):18–42, 2017.
- 302 [12] Kenneth Atz, Francesca Grisoni, and Gisbert Schneider. Geometric deep learning on molecular  
303 representations. *Nature Machine Intelligence*, 3(12):1023–1032, 2021.
- 304 [13] Bowen Jing, Stephan Eismann, Patricia Suriana, Raphael JL Townshend, and Ron Dror. Learn-  
305 ing from protein structure with geometric vector perceptrons. *arXiv preprint arXiv:2009.01411*,  
306 2020.
- 307 [14] Tomohide Masuda, Matthew Ragoza, and David Ryan Koes. Generating 3d molecular struc-  
308 tures conditional on a receptor binding site with deep generative models. *arXiv preprint*  
309 *arXiv:2010.14442*, 2020.
- 310 [15] Shitong Luo, Jiaqi Guan, Jianzhu Ma, and Jian Peng. A 3d generative model for structure-based  
311 drug design. *Advances in Neural Information Processing Systems*, 34:6229–6239, 2021.
- 312 [16] Xingang Peng, Shitong Luo, Jiaqi Guan, Qi Xie, Jian Peng, and Jianzhu Ma. Pocket2mol:  
313 Efficient molecular sampling based on 3d protein pockets. In *International Conference on*  
314 *Machine Learning*, pages 17644–17655. PMLR, 2022.
- 315 [17] Jiaqi Guan, Wesley Wei Qian, Xingang Peng, Yufeng Su, Jian Peng, and Jianzhu Ma. 3d  
316 equivariant diffusion for target-aware molecule generation and affinity prediction. *arXiv preprint*  
317 *arXiv:2303.03543*, 2023.
- 318 [18] Arne Schneuing, Yuanqi Du, Charles Harris, Arian Jamasb, Ilia Igashov, Weitao Du, Tom  
319 Blundell, Pietro Lió, Carla Gomes, Max Welling, et al. Structure-based drug design with  
320 equivariant diffusion models. *arXiv preprint arXiv:2210.13695*, 2022.
- 321 [19] Benoit Baillif, Jason Cole, Patrick McCabe, and Andreas Bender. Deep generative models for  
322 3d molecular structure. *Current Opinion in Structural Biology*, 80:102566, 2023.
- 323 [20] G Richard Bickerton, Gaia V Paolini, Jérémy Besnard, Sorel Muresan, and Andrew L Hopkins.  
324 Quantifying the chemical beauty of drugs. *Nature chemistry*, 4(2):90–98, 2012.
- 325 [21] Christopher A Lipinski, Franco Lombardo, Beryl W Dominy, and Paul J Feeney. Experimental  
326 and computational approaches to estimate solubility and permeability in drug discovery and  
327 development settings. *Advanced drug delivery reviews*, 64:4–17, 2012.
- 328 [22] Zaixi Zhang, Yaosen Min, Shuxin Zheng, and Qi Liu. Molecule generation for target pro-  
329 tein binding with structural motifs. In *The Eleventh International Conference on Learning*  
330 *Representations*, 2022.
- 331 [23] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*  
332 *arXiv:1312.6114*, 2013.
- 333 [24] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
334 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *arXiv preprint*  
335 *arXiv:1406.2661*, 2014.
- 336 [25] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*  
337 *in Neural Information Processing Systems*, 33:6840–6851, 2020.
- 338 [26] Haitao Lin, Yufei Huang, Meng Liu, Xuanjing Li, Shuiwang Ji, and Stan Z Li. Diffbp: Gener-  
339 ative diffusion of 3d molecules for target protein binding. *arXiv preprint arXiv:2211.11214*,  
340 2022.
- 341 [27] Jiaqi Guan, Xiangxin Zhou, Yuwei Yang, Yu Bao, Jian Peng, Jianzhu Ma, Qiang Liu, Liang  
342 Wang, and Quanquan Gu. Decompdiff: Diffusion models with decomposed priors for structure-  
343 based drug design. 2023.
- 344 [28] Meng Liu, Youzhi Luo, Kanji Uchino, Koji Maruhashi, and Shuiwang Ji. Generating 3d  
345 molecules for target protein binding. *arXiv preprint arXiv:2204.09410*, 2022.

- 346 [29] Martin Butterschoen, Garrett M Morris, and Charlotte M Deane. Posebusters: Ai-based docking  
347 methods fail to generate physically valid poses or generalise to novel sequences. *arXiv preprint*  
348 *arXiv:2308.05777*, 2023.
- 349 [30] Cédric Bouysset and Sébastien Fiorucci. Prolif: a library to encode molecular interactions as  
350 fingerprints. *Journal of Cheminformatics*, 13:1–9, 2021.
- 351 [31] Gilles Marcou and Didier Rognan. Optimizing fragment and scaffold docking by use of  
352 molecular interaction fingerprints. *Journal of chemical information and modeling*, 47(1):  
353 195–207, 2007.
- 354 [32] Srinivas Ramachandran, Pradeep Kota, Feng Ding, and Nikolay V Dokholyan. Automated  
355 minimization of steric clashes in protein structures. *Proteins: Structure, Function, and Bioinfor-*  
356 *matics*, 79(1):261–270, 2011.
- 357 [33] Rosa Buonfiglio, Maurizio Recanatini, and Matteo Masetti. Protein flexibility in drug discovery:  
358 from theory to computation. *ChemMedChem*, 10(7):1141–1148, 2015.
- 359 [34] Emanuele Perola and Paul S Charifson. Conformational analysis of drug-like molecules bound  
360 to proteins: an extensive study of ligand reorganization upon binding. *Journal of medicinal*  
361 *chemistry*, 47(10):2499–2510, 2004.
- 362 [35] Anthony K Rappé, Carla J Casewit, KS Colwell, William A Goddard III, and W Mason  
363 Skiff. Uff, a full periodic table force field for molecular mechanics and molecular dynamics  
364 simulations. *Journal of the American chemical society*, 114(25):10024–10035, 1992.
- 365 [36] David Ryan Koes, Matthew P Baumgartner, and Carlos J Camacho. Lessons learned in empirical  
366 scoring with smina from the csar 2011 benchmarking exercise. *Journal of chemical information*  
367 *and modeling*, 53(8):1893–1904, 2013.
- 368 [37] Paul G Francoeur, Tomohide Masuda, Jocelyn Sunseri, Andrew Jia, Richard B Iovanisci, Ian  
369 Snyder, and David R Koes. Three-dimensional convolutional neural networks and a cross-  
370 docked data set for structure-based drug design. *Journal of chemical information and modeling*,  
371 60(9):4200–4215, 2020.
- 372 [38] George C Pimentel and AL McClellan. Hydrogen bonding. *Annual Review of Physical*  
373 *Chemistry*, 22(1):347–385, 1971.
- 374 [39] Deliang Chen, Numan Oezguen, Petri Urvil, Colin Ferguson, Sara M Dann, and Tor C Savidge.  
375 Regulation of protein-ligand binding affinity by hydrogen bond pairing. *Science advances*, 2(3):  
376 e1501240, 2016.
- 377 [40] ID Brown. On the geometry of o–h–o hydrogen bonds. *Acta Crystallographica Section A:*  
378 *Crystal Physics, Diffraction, Theoretical and General Crystallography*, 32(1):24–31, 1976.
- 379 [41] Ylva Andersson, Erika Hult, Henrik Rydberg, Peter Apell, Bengt I Lundqvist, and David C  
380 Langreth. Van der waals interactions in density functional theory. *Electronic Density Functional*  
381 *Theory: Recent Progress and New Directions*, pages 243–260, 1998.
- 382 [42] Elizabeth Barratt, Richard J Bingham, Daniel J Warner, Charles A Laughton, Simon EV Phillips,  
383 and Steve W Homans. Van der waals interactions dominate ligand- protein association in a  
384 protein binding site occluded from solvent water. *Journal of the American Chemical Society*,  
385 127(33):11827–11834, 2005.
- 386 [43] Emily E Meyer, Kenneth J Rosenberg, and Jacob Israelachvili. Recent progress in understanding  
387 hydrophobic interactions. *Proceedings of the National Academy of Sciences*, 103(43):15739–  
388 15746, 2006.
- 389 [44] Rohan Patil, Suranjana Das, Ashley Stanley, Lumbani Yadav, Akulapalli Sudhakar, and Ashok K  
390 Varma. Optimized hydrophobic interactions and hydrogen bonding at the target-ligand interface  
391 leads the pathways of drug-designing. *PloS one*, 5(8):e12029, 2010.
- 392 [45] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ron-  
393 neberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al.  
394 Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.

- 395 [46] Andrew M Davis and Simon J Teague. Hydrogen bonding, hydrophobic interactions, and failure  
396 of the rigid receptor hypothesis. *Angewandte Chemie International Edition*, 38(6):736–749,  
397 1999.
- 398 [47] Ting Chen. On the importance of noise scheduling for diffusion models. *arXiv preprint*  
399 *arXiv:2301.10972*, 2023.
- 400 [48] Wm H Rinkenbach. The heats of combustion and formation of aromatic nitro compounds.  
401 *Journal of the American Chemical Society*, 52(1):115–120, 1930.
- 402 [49] Bernard R Brooks, Charles L Brooks III, Alexander D Mackerell Jr, Lennart Nilsson, Robert J  
403 Petrella, Benoît Roux, Youngdo Won, Georgios Archontis, Christian Bartels, Stefan Boresch,  
404 et al. Charmm: the biomolecular simulation program. *Journal of computational chemistry*, 30  
405 (10):1545–1614, 2009.
- 406 [50] Irina Kufareva, Andrey V Ilatovskiy, and Ruben Abagyan. Pocketome: an encyclopedia of  
407 small-molecule binding sites in 4d. *Nucleic acids research*, 40(D1):D535–D540, 2012.
- 408 [51] Martin Steinegger and Johannes Söding. Mmseqs2 enables sensitive protein sequence searching  
409 for the analysis of massive data sets. *Nature biotechnology*, 35(11):1026–1028, 2017.
- 410 [52] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics.  
411 In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages  
412 681–688, 2011.
- 413 [53] John Ingraham, Max Baranov, Zak Costello, Vincent Frappier, Ahmed Ismail, Shan Tie, Wujie  
414 Wang, Vincent Xue, Fritz Obermeyer, Andrew Beam, et al. Illuminating protein space with a  
415 programmable generative model. *bioRxiv*, pages 2022–12, 2022.
- 416 [54] Jason Yim, Brian L Trippe, Valentin De Bortoli, Emile Mathieu, Arnaud Doucet, Regina  
417 Barzilay, and Tommi Jaakkola. Se (3) diffusion model with application to protein backbone  
418 generation. *arXiv preprint arXiv:2302.02277*, 2023.

## 419 Checklist

- 420 1. For all authors...
- 421 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
422 contributions and scope? [Yes]
- 423 (b) Did you describe the limitations of your work? [Yes]
- 424 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 425 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
426 them? [Yes]
- 427 2. If you are including theoretical results...
- 428 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 429 (b) Did you include complete proofs of all theoretical results? [N/A]
- 430 3. If you ran experiments (e.g. for benchmarks)...
- 431 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
432 mental results (either in the supplemental material or as a URL)? [Yes]
- 433 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
434 were chosen)? [Yes]
- 435 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
436 ments multiple times)? [Yes]
- 437 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
438 of GPUs, internal cluster, or cloud provider)? [N/A]
- 439 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 440 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 441 (b) Did you mention the license of the assets? [Yes] In the repository

- 442 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
443 (d) Did you discuss whether and how consent was obtained from people whose data you're  
444 using/curating? [N/A] Publicly available.  
445 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
446 information or offensive content? [N/A]  
447 5. If you used crowdsourcing or conducted research with human subjects...  
448 (a) Did you include the full text of instructions given to participants and screenshots, if  
449 applicable? [N/A]  
450 (b) Did you describe any potential participant risks, with links to Institutional Review  
451 Board (IRB) approvals, if applicable? [N/A]  
452 (c) Did you include the estimated hourly wage paid to participants and the total amount  
453 spent on participant compensation? [N/A]

## 454 A CrossDocked Dataset

455  
456 The CrossDocked dataset is a standard dataset used in the field of generative modelling for structure-  
457 based drug design [37]; since the models benchmarked here were trained on this dataset, it is the  
458 benchmarking dataset of choice. It was originally created by clustering PDB structures by "pocket  
459 similarity" via Pocketome [50], i.e. grouping structures with similar ligand binding sites together.  
460 To expand the dataset beyond this initial data, all ligands with a molecular weight < 1000 Da that  
461 were associated with a given pocket were docked into each receptor assigned to that pocket via  
462 the docking tool smina [36]. This cross-docking process results in the basis dataset CrossDocked  
463 2020 [37], which contains 2,922 pockets, 18,450 complexes and 13,839 ligands, together comprising  
464 around 22.5 million poses (i.e. protein-ligand structures).

465 Most generative models are however not trained on this raw dataset, but on a filtered version of it,  
466 following the procedure of the Pocket2Mol model [16]. As a quality control, data points whose  
467 binding pose RMSD is greater than 1 were filtered out. This leads to a filtered dataset with 184,057  
468 data points. The mmseq2 program [51] was used to cluster data at 30% identity, and training and test  
469 sets were created by randomly drawing 100,000 protein-ligand pairs for training and 100 proteins  
470 from the remaining clusters for testing.

471 The 100 proteins comprising the test set are on average around 320 residues long, with the biggest  
472 protein having a length of 752 residues.

## 473 B Extended Implementation

474

### 475 B.1 Methods Implementation

476 All generative methods accessed were trained using the same dataset and splits as proposed in  
477 Peng et al. [16]. Docking protocols were done using the SMINA settings described in the original  
478 CrossDocked paper [37].

479

### 480 B.2 Procedure of model reproduction

481 For generated poses, we sourced molecules from Schneuing et al. [18] for DiffSBDD, and Guan et al.  
482 [17] for CrossDocked, TargetDiff, Pocket2Mol, 3DSBDD and LiGAN (where they provide generated  
483 poses but we additionally perform our own redocking).

484 For FLAG [22], no weights were provided so we retrained the model as described in Zhang et al.  
485 [22] using the code and config file available at [github.com/zaixizhang/FLAG](https://github.com/zaixizhang/FLAG). When sampling,  
486 we found that generation was attempted 100 times per target and then any molecules with fewer than  
487 8 atoms were discarded. This ended up encompassing the majority of molecules, resulting in small  
488 test sizes, so we implemented a while loop to sample 100 molecules whilst keeping faithful to the  
489 filtering used in the codebase. Having modified the code to work on GPU, sampling 100 targets took  
490 about 1-2 minutes per target on a single A100 GPU.

491 For DecompDiff [27], we use the official implementation with the published weights available at  
492 [github.com/bytedance/DecompDiff](https://github.com/bytedance/DecompDiff). We sampled 100 samples for each of the 100 targets using  
493 the `sample_diffusion_drift.py` script in `ref_prior` mode. With the provided code, sampling  
494 100 targets took about 20-30 minutes per target on a single A100 GPU.

## 495 C Recommendations for future work

496 **Exploring reduced-noise sampling strategies** Interestingly, both diffusion-based works (DiffS-  
497 BDD and TargetDiff) performed similarly in terms of strain energy (see Section 4.5). We hypothesize  
498 this may be due to the injection of random noise into the coordinate features at all but the last step of  
499 stochastic gradient Langevin dynamics samplings [52], making it challenging to construct precise  
500 bond angles etc. Here, inspiration could be taken from protein design. For example, Chroma develops  
501 a low-temperature sampling regime to reduce the effect of noise [53], FrameDiff effectively scales

502 down injected noise [54], both resulting in a substantial increase in sample quality with an acceptable  
503 decrease in sample diversity.

504 **Heavily penalise steric clashes during training** All evaluated methods frequently create steric  
505 clashes, resulting in physically unrealizable samples. We suggest that mitigating steric clashes is  
506 key for the next generation of SBDD models. This could be done via extra loss terms, for example,  
507 by including a distogram loss as in AlphaFold2 [45] or the steric clash loss in LiGAN [14] and  
508 DecompDiff [27] (note that later method does not explicitly measure clashes). A similar loss-based  
509 approach has been effective in mitigating chain-breaks diffusion models for protein backbone design  
510 [54].

511 **Consider representing hydrogens** Virtually all work in ML for structural biology chooses to  
512 not explicitly represent hydrogen atoms [45, 16, 15, 54, 18, 17], under the assumption that they  
513 can be *implicitly* learned and reasoned over with deep neural networks. However, our analysis of  
514 hydrogen bond networks within generated molecules found that generative methods struggle to handle  
515 the precise geometries required to make a hydrogen bond [40] (even when redocked). Despite the  
516 increased computational cost, we therefore recommend that future work explores their inclusion.

## 517 D Additional Figures

### 518 D.1 Interactions analysis

519 We include the comparisons between generative method against baselines for both Van der Waals  
520 contacts and hydrophobic interactions, both for generated redocked poses in Figure 6.

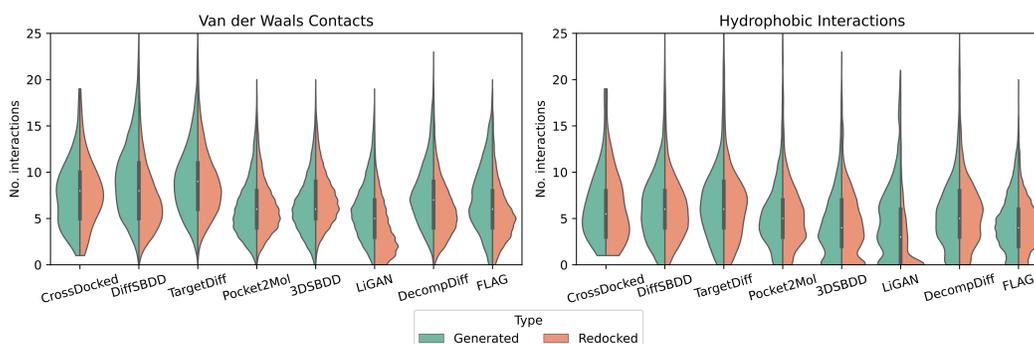


Figure 6: Extended analysis of the interaction profiles of the generated molecules for the different methods. While the focus in the main text was on hydrogen bonds, the results in this figure include Van der Waals Contacts and hydrophobic interactions, reported for both the generated as well as the redocked pose.

### 521 D.2 Redocking and clashes analysis

522 In Figure 7, we provide the per target redocking RMSDs per method. Figure 8 and 9 show the number  
523 of steric clashes per target for the generated and redocked poses respectively.

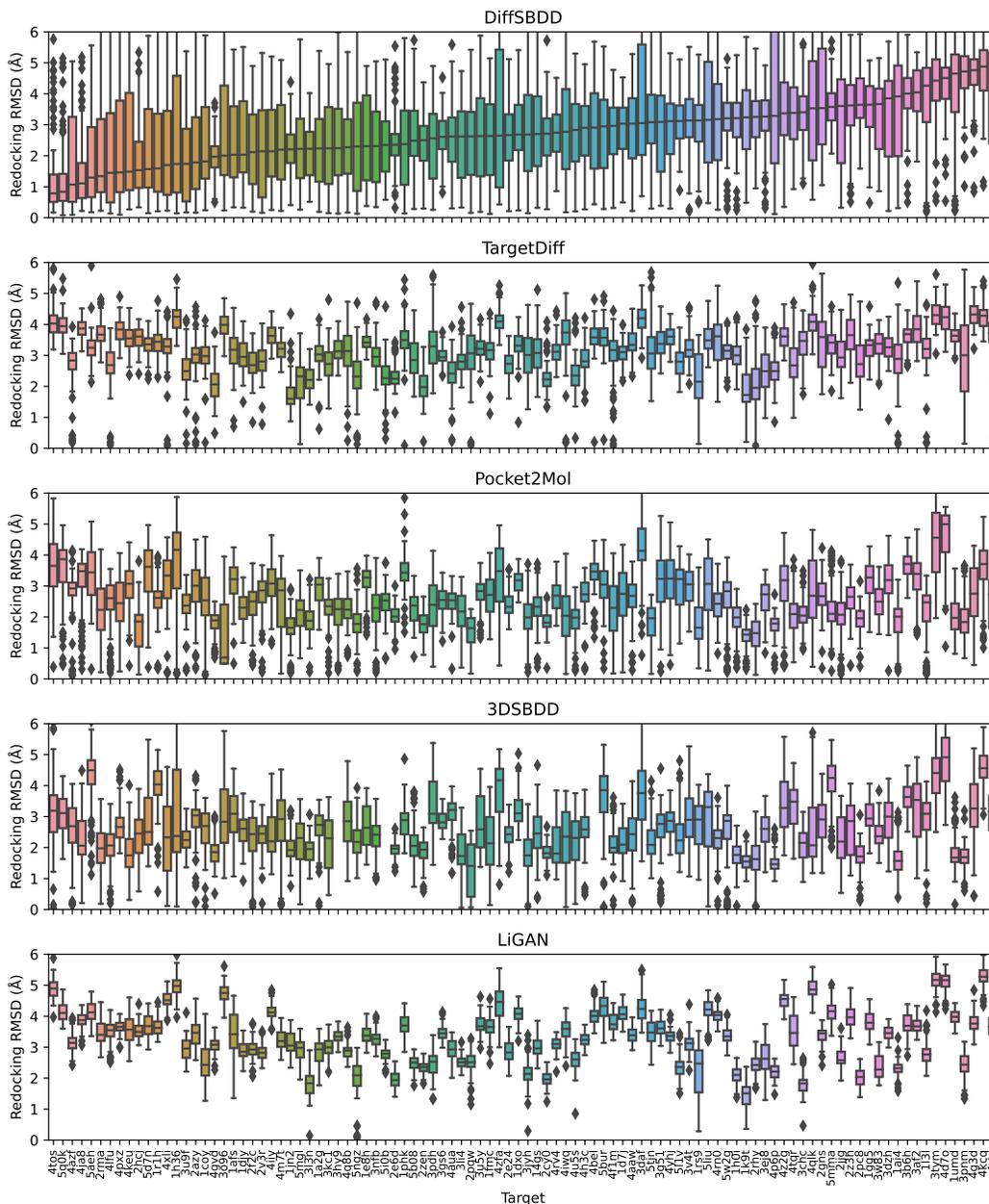


Figure 7: Redocking RMSD per method per target for CrossDocked test set. Order is determined arbitrarily by median score per target for DiffSBDD.

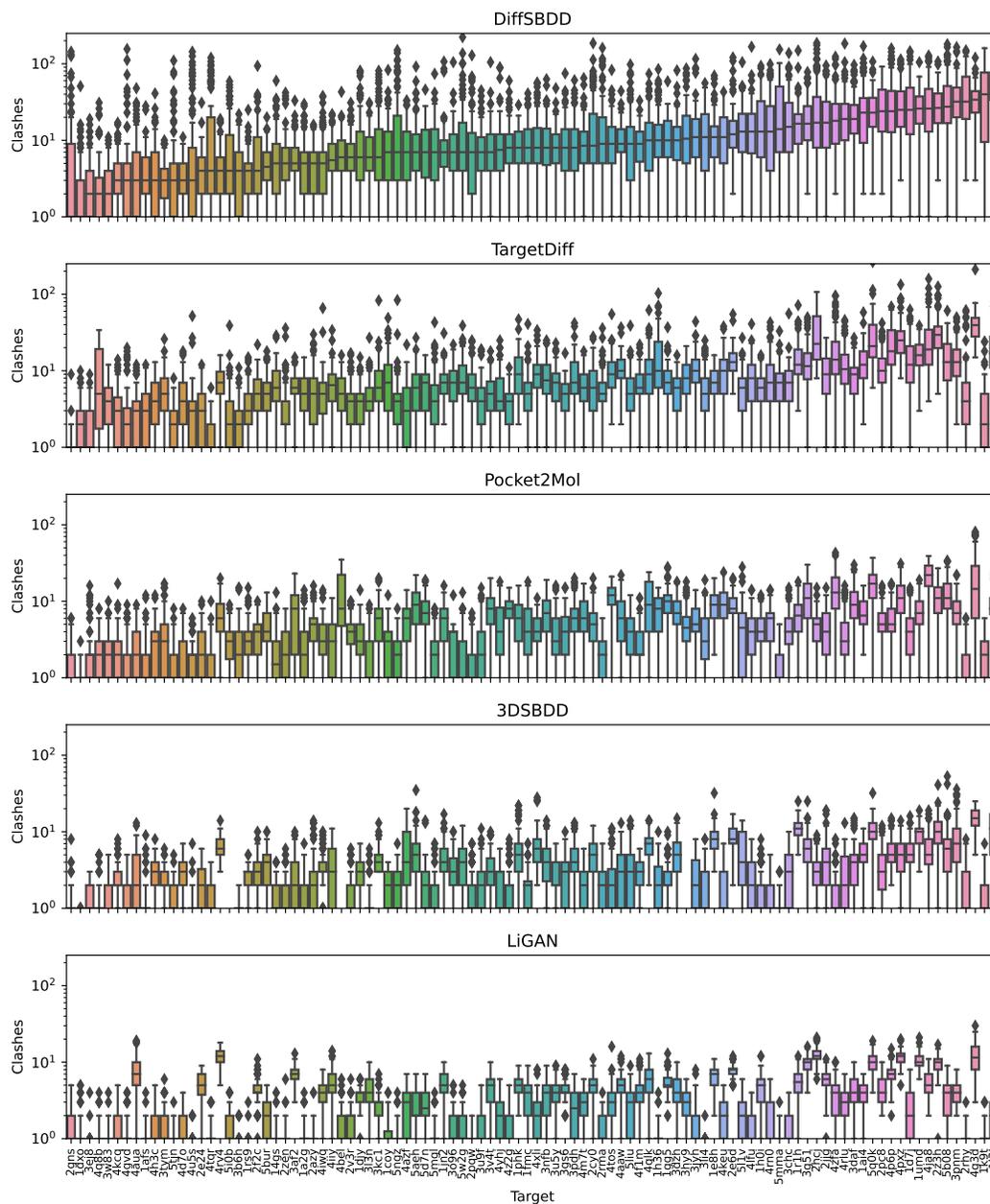


Figure 8: Steric clashes per method per target for generated poses in the CrossDocked test set. Order is determined arbitrarily by median score per target for DiffSBDD.

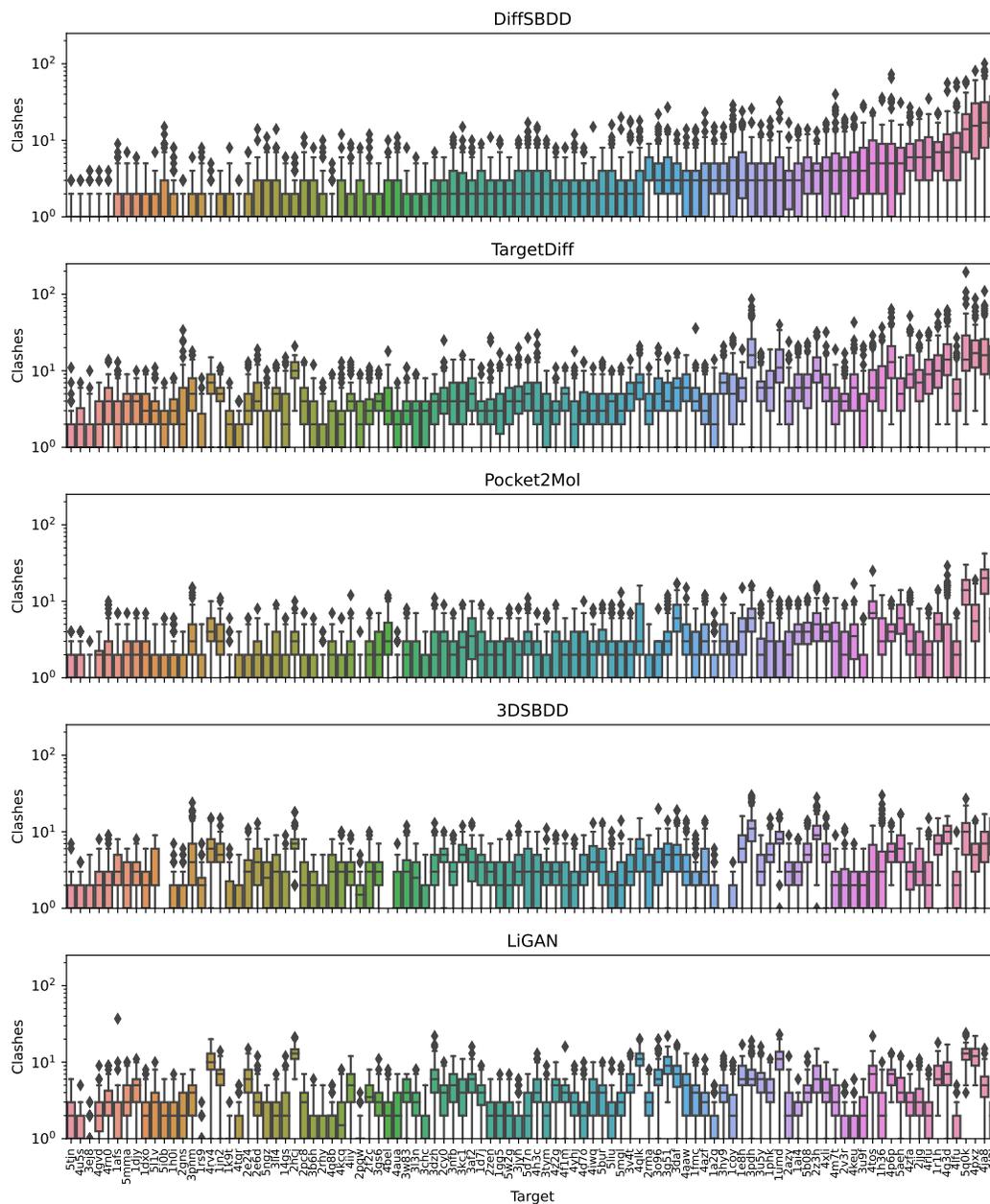


Figure 9: Steric clashes per method per target for redocked poses in the CrossDocked test set. Order is determined arbitrarily by median score per target for DiffSBDD.