

MoleGuLAR: Molecule Generation Using Reinforcement Learning with Alternating Rewards

Manan Goel, Shampa Raghunathan, Siddhartha Laghuvarapu, and U. Deva Priyakumar*



Cite This: <https://doi.org/10.1021/acs.jcim.1c01341>



Read Online

ACCESS |



Metrics & More

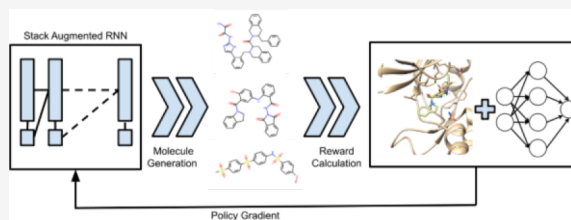


Article Recommendations



Supporting Information

ABSTRACT: The design of new inhibitors for novel targets is a very important problem especially in the current scenario with the world being plagued by COVID-19. Conventional approaches such as high-throughput virtual screening require extensive combing through existing data sets in the hope of finding possible matches. In this study, we propose a computational strategy for de novo generation of molecules with high binding affinities to the specified target and other desirable properties for druglike molecules using reinforcement learning. A deep generative model built using a stack-augmented recurrent neural network initially trained to generate druglike molecules is optimized using reinforcement learning to start generating molecules with desirable properties like LogP, quantitative estimate of drug likeliness, topological polar surface area, and hydration free energy along with the binding affinity. For multiobjective optimization, we have devised a novel strategy in which the property being used to calculate the reward is changed periodically. In comparison to the conventional approach of taking a weighted sum of all rewards, this strategy shows an enhanced ability to generate a significantly higher number of molecules with desirable properties.



1. INTRODUCTION

The advent of data-driven techniques across multiple domains of computer science, such as robotics, natural language processing, and computer vision, has found immense success, and this has led to their application in natural sciences.^{1,2} The curation of large data sets^{3–5} has increased the relevance of machine-learning-based approaches in problems like molecular property prediction, conceiving retrosynthetic pathways, protein structure prediction, and drug discovery.^{6–9}

Drug discovery is a long, expensive, and arduous process which combines a wide range of disciplines including chemistry, biology, and pharmacology. For a novel target, the conventional approach is to perform high-throughput screening on chemical libraries to identify small molecules that bind well to the target. The identified hits are then optimized to get higher binding affinity, reduce toxicity, and improve oral bioavailability.^{10,11} The time and expense involved in this process give rise to alternate in silico approaches like virtual screening wherein small molecules from existing drug libraries are computationally evaluated by generating protein ligand complexes using docking calculations and ranking them using a scoring function.^{12,13} However, these also come with the caveat that finding the most stable conformation of the complex is a nonconvex optimization problem, and it can take a very large amount of time (≈ 10 min per molecule) to find the most optimal conformation. These can be made faster using machine-learning-based approaches like the works of Aggarwal et al.¹⁴ for detecting the ligand-binding site, Chelur and Priyakumar¹⁵ for binding residue detection, and Mehta et al.¹⁶ for enhanced molecular sampling.^{14–16} However, even the most exhaustive studies¹⁷ have been able to find

binding affinities of $\approx 10^8$ molecules on a single target which is minuscule in comparison to the vast magnitude of the chemical space with about 10^{60} synthesizable molecules.¹⁸ This posits the argument for the de novo generation of molecules with high binding affinities to the required target instead of searching in existing libraries.

Machine-learning-based approaches like recurrent neural networks, generative adversarial networks (GANs), and variational autoencoders have recently been adopted for molecule generation. Gupta et al.¹⁹ used long short-term memory recurrent neural networks, generally used for natural language processing tasks, to generate molecules in the form of SMILES (simplified molecular-input line-entry system), which is a string representation of molecules and has its own grammar and semantics.²⁰ GANs are generative models that learn the probability distribution of the training data, and sampling from the distribution can then be used to generate synthetic data points. This model has also been applied to the generation of molecules with desirable properties in works by De Cao and Kipf,²¹ Prykhodko et al.,²² and Maziarka et al.²³ Jin et al.²⁴ used the graph representation of molecules to train a variational autoencoder that could then generate graphs of new

Received: November 5, 2021

69 molecules.²⁴ Kusner et al.,²⁵ Griffiths and Hernández-Lobato,²⁶
 70 and Lim et al.²⁷ used SMILES representations for generating
 71 molecules through the variational autoencoder architec-
 72 ture.^{25–27} In fact, applications of deep learning models for
 73 molecule generation have proven to be very successful in recent
 74 years.^{19,28–30}

75 The next challenge is to generate molecules with desirable
 76 properties for which the two major approaches being adopted
 77 are reinforcement learning and latent space optimization.
 78 Variational autoencoders are capable of learning a continuous
 79 space representation of molecules^{31–34} which can then be
 80 optimized to generate molecules with target properties through
 81 techniques like Bayesian optimization and swarm optimiza-
 82 tion.^{35,36} Gao et al.³⁷ used an autoencoder architecture and a
 83 generator model in combination to incrementally update the
 84 latent space representation of the given molecule to reach a
 85 molecule with the desired properties.³⁷ Reinforcement learning
 86 can be used to generate desirable molecules by decomposing the
 87 process as a sequence of states and actions to maximize a reward
 88 which in this case is the desirable property. Popova and co-
 89 workers used stack-augmented gated recurrent units (GRUs) to
 90 generate molecules followed by reinforcement learning guided
 91 optimization on properties like LogP, quantitative estimate of
 92 drug likeliness (QED), and synthetic accessibility.³⁸ Guimaraes
 93 et al.³⁹ combined the GAN and reinforcement learning
 94 frameworks for the task while You et al.⁴⁰ and Khemchandani
 95 et al.⁴¹ used a graph-based policy network to generate molecular
 96 graphs.^{39–41} Boitreaud et al.⁴² and Born et al.⁴³ combined a
 97 variational autoencoder model with reinforcement learning to
 98 generate molecules with high binding affinities to the specified
 99 target and antiviral candidates, respectively.^{42–44} Bung et al.⁴⁵
 100 used reinforcement learning with a SMILES generator to
 101 generate molecules with high binding affinities to JAK2 and
 102 SARS-CoV-2 3CL proteins, respectively.^{45,46}

103 In this work we propose a molecule generation pipeline,
 104 MoleGuLAR (molecule generation using reinforcement learn-
 105 ing with alternating rewards, Figure 1), which uses a stack
 106 augmented recurrent neural network (RNN) initially trained to
 107 generate valid druglike molecules which is then optimized to
 108 generate molecules with a high binding affinity to the specified
 109 target. For the binding affinity calculation, we tried two
 110 methodologies: (1) performing docking calculation to find the

111 most stable complex and the corresponding binding affinity and
 112 (2) using a machine-learning model trained to predict binding
 113 affinities. In the case of multiobjective optimization, we found
 114 that using a weighted sum of the rewards from different
 115 properties may not be effective in some cases because it is
 116 possible that one or more properties dominate others leading to
 117 poor results. Hence, we also propose a novel optimization
 118 strategy in which the reward is alternated so that the model
 119 changes the region from which it samples in the chemical space.
 120 When the reward is changed, the generator starts from a better
 121 position with respect to one property when optimization for
 122 another property is started. We also showcase its application on
 123 two proteins: M_{pro} of SARS-CoV-2 and TBK1 with a wide set
 124 of target properties. The robustness of this strategy is further
 125 showcased by using it to optimize the model for conflicting
 126 properties along with the binding affinity.

2. THEORY AND METHODS

127 This section describes the various components of the proposed
 128 framework (Figure 1). The formulation of the stack-augmented
 129 RNN as the generator model is detailed in Section 2.1 followed
 130 by methods for binding affinity calculation and hydration free
 131 energy calculation in Sections 2.2 and 2.3, respectively. The
 132 formulation of the molecule generation as a Markov decision
 133 process, use of reinforcement learning to maximize a given
 134 reward function using policy gradient, and the two optimization
 135 strategies used in this study are described in Section 2.4.

136 **2.1. Generator.** The generator module makes use of a stack-
 137 augmented GRU which outputs molecules as SMILES strings as
 138 presented by Popova et al.^{38,47} A valid SMILES string must have
 139 correct valency for all atoms, and all ring openings and closures
 140 must be counted; hence, conventional RNNs do not work well
 141 on this task because of their inability to count. Therefore, the
 142 addition of a memory unit along with the RNN forms an
 143 appropriate model which is explained further in Section S1.1 of
 144 the Supporting Information.

145 The stack RNN is initially trained on ≈ 1.5 million druglike
 146 molecules from the ChEMBL21 database⁵ to learn the rules and
 147 grammar of SMILES strings.

148 **2.2. Binding Affinity Calculation.** **2.2.1. Docking Calculations.**
 149 The generator model once trained is then used to
 150 produce ligands which are then docked to the specified target to
 151 find the most stable conformation of the complex and to find the
 152 corresponding binding affinity further referred to as BA in the
 153 article. The 3D structure of the molecule from the SMILES
 154 string is obtained using the RDKit toolkit.⁴⁸ The target proteins
 155 TBK1 (PDB ID: 4BTK) and SARS-CoV-2 M_{pro} complexed
 156 with N3 inhibitor (PDB ID: 6LU7) are obtained from the
 157 Research Collaboratory for Structural Bioinformatics PDB.⁴⁹
 158 The ligand and protein structure is then converted to a format
 159 suitable for the input to the docking software using
 160 AutoDockTools4. The molecule docking grid is generated in
 161 the next step using the AutoGrid4 utility, and finally the docking
 162 calculation is done using AutoDock-GPU while keeping the
 163 protein active site rigid.^{50,51} This tool is referred as AutoDock
 164 in the rest of the article. Detailed information about the docking
 165 methodology is provided in the Supporting Information.

166 **2.2.2. Machine-Learning Models.** We also tested the use of
 167 machine-learning models as a placeholder for AutoDock to
 168 predict the binding affinities of the generated ligands instead of
 169 performing docking calculations to reduce the computation
 170 time. In order to do this, we obtained a data set of ≈ 2 million

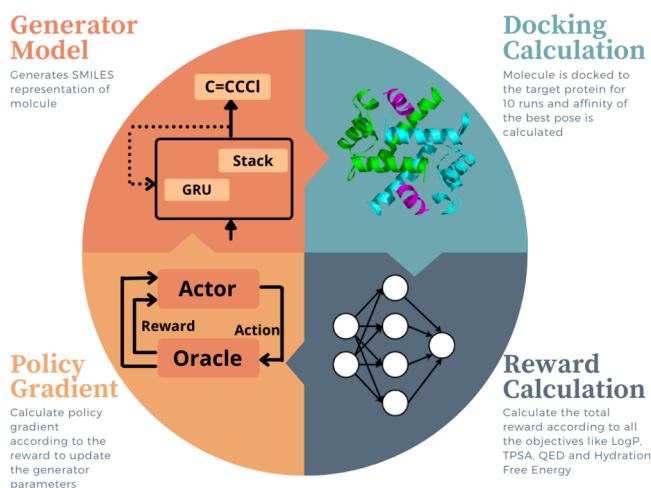


Figure 1. Pipeline used by MoleGuLAR for generating molecules with high binding affinity to a specified drug target.

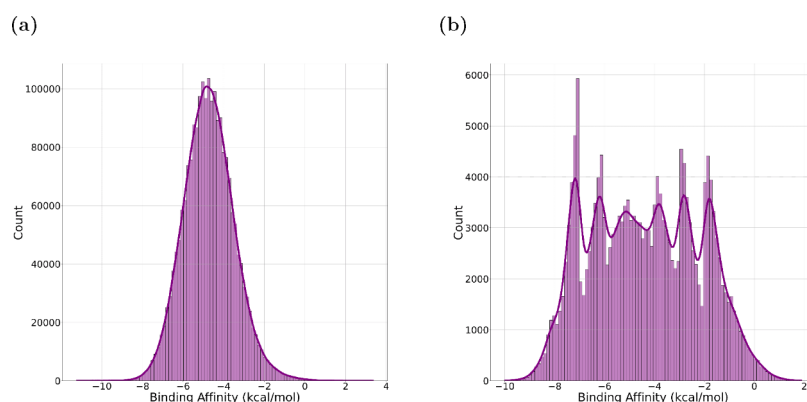


Figure 2. Distribution of binding affinities. (a) Two million molecules with TTBK1 and (b) selected molecules from buckets.

171 molecules obtained from the HTS collection by Enamine⁵²
172 docked with the TTBK1 protein.

173 **Figure 2a** shows that a significant number of molecules lie in a
174 small range of binding affinities, and hence, using that for the
175 predictor model tends to overfit (Figure S2 of the [Supporting](#)
176 [Information](#)). In order to tackle this issue, we split the entire data
177 set into smaller bins of 1 kcal/mol and sampled 25K molecules
178 from each bin and all the molecules if the number of molecules is
179 less than 25K. **Figure 2b** shows the distribution of the obtained
180 subset consisting of ≈ 200 K molecules. This is then split into
181 training, testing, and validation sets in the ratio 80:10:10.

182 We further tested various machine-learning models for this
183 regression task. Jaeger et al.⁵³ proposed the Mol2vec⁵³ model for
184 learning vector representations of SMILES strings that can then
185 be used as input for further downstream tasks like binding
186 affinity prediction as done by Mehta et al.¹⁶ along with
187 predicting other properties. Using these embeddings, a random
188 forest model with 250 decision trees was trained for predicting
189 the BA. The aforementioned model with input features from the
190 embeddings obtained from graph isomorphism networks (GIN)
191 proposed by Xu et al.⁵⁴ and Hu et al.⁵⁵ were also used for the
192 task. The drawback of these approaches is that the embeddings
193 being used remain constant during training and that leads to
194 poor performance ([Figure S3](#) of the [Supporting Information](#)).
195 Fine tuning these representations during the training process
196 helps to improve the accuracy for which three linear layers were
197 added with the GIN embeddings, and the model is then trained
198 end-to-end.

199 **2.3. Hydration Free Energy Prediction.** The hydration
200 free energy (ΔG_{Hyd}) of a molecule measures its interaction with
201 water and forms an important part of the drug delivery system.
202 The state of the art methods for predicting it include message-
203 passing neural networks (MPNN)⁹ as shown by Wu et al.³ in
204 MoleculeNet³ and chemically interpretable graph interaction
205 networks (CIGIN)⁵⁶ by Pathak et al.⁵⁶ To predict ΔG_{Hyd} ,
206 presently we took out 10% molecules out of 643 molecules as a
207 hold out test set and performed fivefold cross-validation on the
208 remaining data set with 10% going in the validation set and 80%
209 in the training set.

210 **2.4. Reinforcement Learning.** A reinforcement-learning
211 pipeline generally consists of two modules: the actor and the
212 critic. The actor takes the current state (s_t) of the system and
213 performs an action (a_t) that should maximize the reward. The
214 critic sees a_t , s_t , and the state obtained by performing the action
215 (s_{t+1}) and penalizes or rewards the actor.

216 Generation of a SMILES string can be modeled as a Markov
217 decision process where s_t denotes the SMILES string

constructed so far, and a_t denotes the addition of a token to s_t .
218 We also define a terminal state s_T which signifies the end of the
219 molecule and initial state s_0 . The whole generation process is
220 depicted in [Figure 3](#). 221

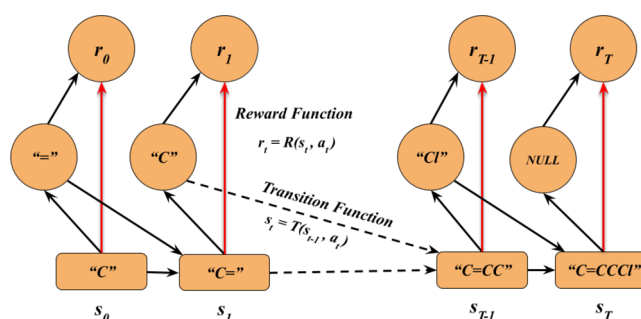


Figure 3. Illustration of SMILES string generation as a Markov decision process. At each state s_t , the agent performs an action a_t to give the updated state s_{t+1} and provide a reward according to the state.

The generator model parametrized by Θ estimates the
222 probability, $p(a_t|s_t, \Theta)$, samples a_t from the probability
223 distribution, and updates the state until s_T is reached. Rewards
224 of all intermediate states s_t with $t < T$ are 0 since it is possible that
225 the intermediate SMILES strings may not represent a valid
226 molecule. s_T is then sent to the critic which returns the reward
227 $r(s_T)$. Hence, the task here is to find Θ such that the expected
228 reward given by [eq 1](#) is maximized. This is done using the
229 REINFORCE algorithm⁵⁷ which is detailed further in the
230 [Supporting Information](#). 231

$$R(\Theta) = \mathbb{E}[r(s_T)|s_0, \Theta] = \sum_{s_T \in \mathcal{S}} p_{\Theta}(s_T) r(s_T) \quad (1) \quad 232$$

For the multiobjective setup, the reward function $r(s_T)$ is
233 composed of multiple components from the different properties
234 for which the model is being optimized. The two reward
235 strategies that we propose are 236

• **Weighted Sum Rewards:** The total reward $r(s_T)$ is expressed
237 as a weighted sum of all other components: 238

$$r(s_T) = w_1 D(s_T) + w_2 L(s_T) + w_3 Q(s_T) + w_4 T(s_T) \\ + w_5 H(s_T)$$

where D fetches the reward for BA, L for LogP, Q for QED, T for
239 topological polar surface area, and H for ΔG_{Hyd} . Weights are
240 kept as hyperparameters and remain constant throughout the 241

242 optimization process. The functional forms of the reward for
243 each property are given in Table S2 of the Supporting
244 Information.

245 •Alternating Rewards: The aforementioned approach does
246 not work especially in cases where properties are conflicting like
247 high TPSA, and more negative hydration free energy would be
248 contradictory in nature. In such cases we have devised a strategy
249 wherein all the weights are changed to 0 except one. This takes
250 the generator model into the space where one property is
251 optimal providing a better starting point when optimization for
252 the other property is started. The current strategy works
253 extremely well across most of the tasks and removes the
254 requirement for acute hyperparameter tuning to find the most
255 optimal weights for each reward function. Further details are
256 given in the Results and Discussion section.

3. RESULTS AND DISCUSSION

257 This section describes the performance of machine-learning
258 models for predicting binding affinity and ΔG_{Hyd} as well as
259 application of the proposed pipeline on the targets:

Table 1. Performance of Predictor Models for BA in Terms of Performance Metrics MAE (kcal/mol) and Coefficient of Determination (R^2)

model	MAE (kcal/mol)	R^2
graph embeddings + random forest	0.87	0.55
Mol2vec + random forest	0.47	0.91
graph isomorphism network (GIN)	0.45	0.93

- 260 • SARS-CoV-2 M_{pro} (PDB ID: 6LU7): With the world in
261 the midst of a global pandemic caused by COVID-19, the
262 main protease (M_{pro}) has been identified as an important
263 target due its vital role in viral transcription and
264 replication.⁵⁸
- 265 • TTBK1 (PDB ID: 4BTK): Neurodegenerative diseases
266 have become extremely common over the past few years,
267 and the tau-tubulin kinase 1 has proved to be an attractive
268 target to combat a wide variety of neurodegenerative
269 diseases.⁵⁹

270 All of the presented experiments were performed using an
271 Intel Xeon E5-2640 v4 processor and a Nvidia GeForce RTX
272 2080Ti GPU. The implementation details of all the machine-
273 learning models are described in Section S1.3 of the Supporting
274 Information.

3.1. Machine-Learning Predictor Models. **3.1.1. Binding Affinity.** The performance of the random forest model with different embeddings and the performance of GIN discussed in Section 2.2.2 for BA prediction on the test set are reported in Table 1.

The use of constant embeddings for the random forest models leads to a higher mean absolute error (MAE) in comparison to the GIN model because in the latter, the model learns to automatically extract more relevant information from the molecular graph. The former also showed comparatively poorer performance in the desirable region, i.e., where the BA is high due to the lesser number of samples in that range in the data set. The correlation between the predicted values and ground truth values is shown in Figure 4a.

3.1.2. Hydration Free Energy. Figure 4b shows the correlation of predicted and ground truth ΔG_{Hyd} in the test set obtained from the FreeSolv data set. The MPNN model succeeds in achieving a high degree of accuracy with a root mean squared error (RMSE) of 1.35 kcal/mol and close correlation characterized by the R^2 score of 0.87. However, the use of machine-learning models for predicting properties of a molecule comes with the caveat that in the case in which any substructures in the molecule are not present in the training set, the prediction may be inaccurate. This is especially true in regions where the existing data is very sparse like cases when the binding affinity is extremely high. As shown in the subsequent sections, performing docking calculations on the generated molecules addresses this issue.

3.2. Single Objective Optimization. The initial tests were performed to analyze the ability of the generator model based only on SMILES to learn to generate molecules with structure complementary to the binding pocket. In sections 3.2.1 and 3.2.2, we evaluate two different approaches to obtaining the binding affinity, namely docking calculations and using a GIN predictor model, respectively. To further validate the performance, the optimization was done from scratch across three runs with different random seeds, and the distribution of the properties of the generated molecules is given in the Supporting Information. Refer to Section S4 of the Supporting Information for the statistics and generated molecules from each of the trials in Section 3.2.

3.2.1. Docking Calculations. For both TTBK1 and SARS-CoV-2 M_{pro} , the generator model was optimized for 175 iterations with 15 policy gradient steps in each and a batch of 10 molecules. At the end of every iteration, 100 molecules were

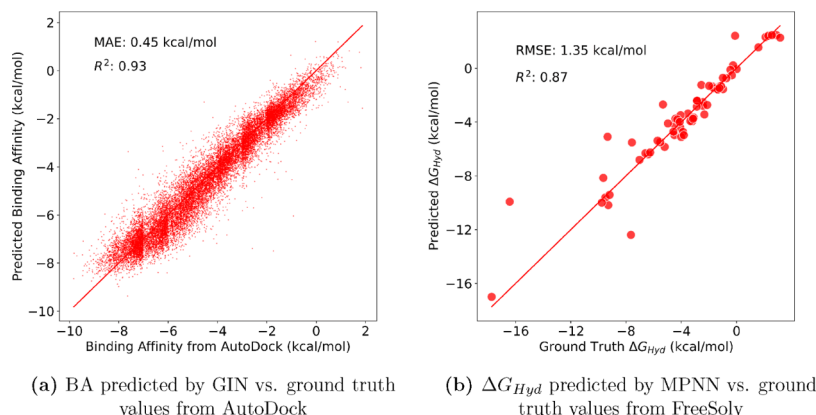


Figure 4. Correlation plots between predicted values from the machine-learning models and ground truth values from the respective data sets.

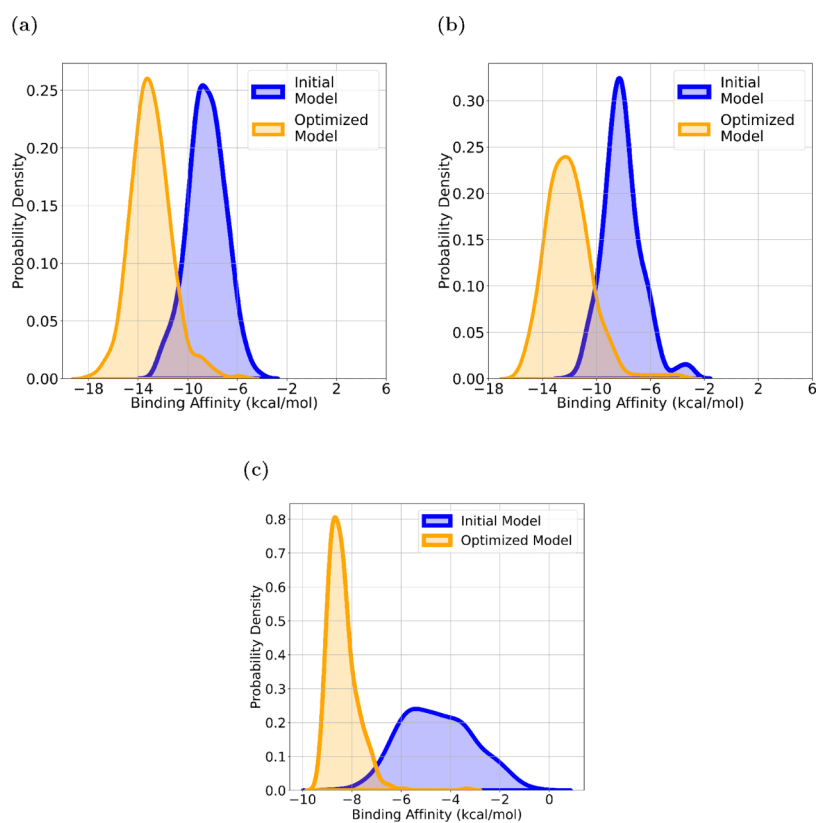


Figure 5. Distribution of BA of generated molecules before and after optimizing the generator for reward from BA with (a) SARS-CoV-2 M_{pro} , (b) TTBK1, and (c) TTBK1 calculated using GIN.

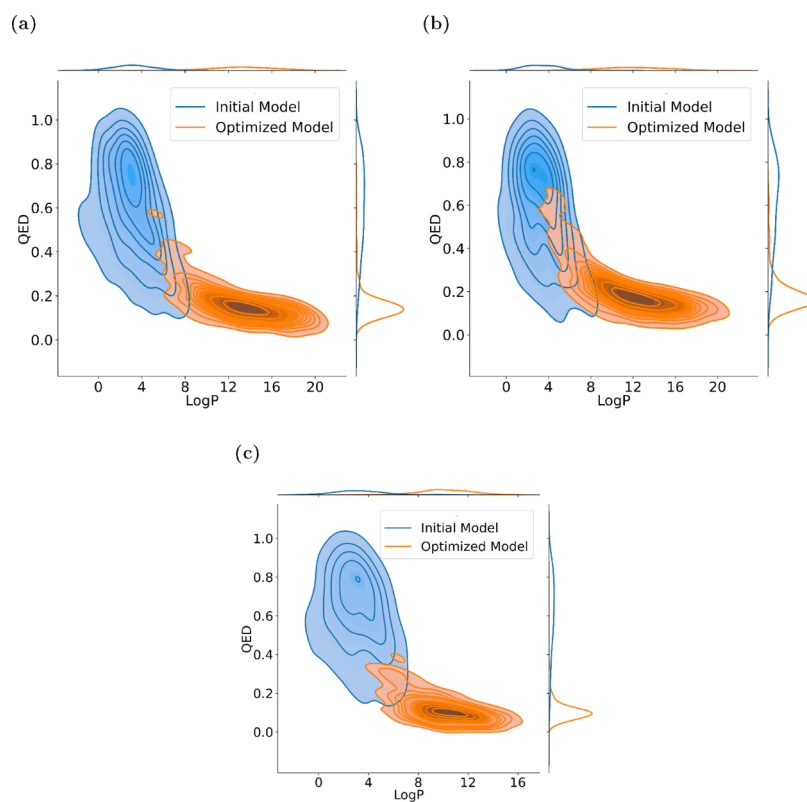


Figure 6. Joint distribution of LogP and QED of generated molecules before and after optimizing the model for reward from BA with (a) SARS-CoV-2 M_{pro} , (b) TTBK1, and (c) TTBK1 calculated using GIN.

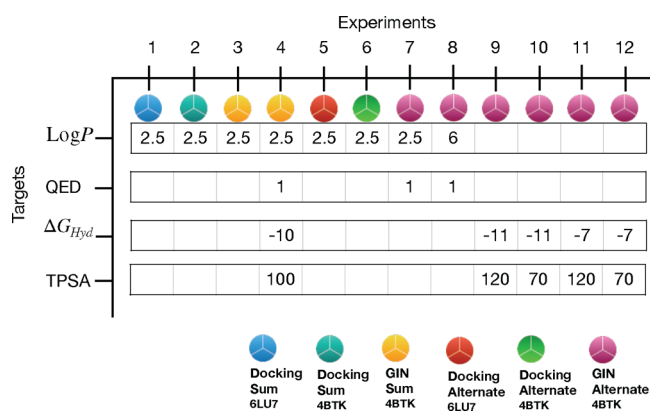


Figure 7. Tests performed for multiobjective optimization: protein PDB ID, tools used for BA calculation, different target values of respective properties, and optimization strategies. ΔG_{Hyd} and TPSA are in kcal/mol and \AA^2 , respectively.

6b). Hence, there is a need for multiobjective optimization in order to keep the other properties in check as well.

3.2.2. *Using GIN.* The use of a GIN for BA prediction leads to an approximately 10-fold speed-up in optimization taking only about 5 h while still keeping the high performance. To further validate the generator, 500 molecules were generated and docked to TTBK1. The shifts in distribution of BA, LogP, and QED are shown in Figures 5c and 6c. However, the undesirable LogP and QED persist, and hence multiobjective optimization is used.

3.3. Multiobjective Optimization. In order to address the shortcomings of single objective optimization, the rewards from different properties were also integrated into the policy gradient calculation, and the performance was tested for the two proteins using different target values and both optimization strategies. These are listed in Figure 7, and the results have been discussed in the subsequent sections. Section 3.3.1 describes the conventional method for multiobjective optimization by taking a weighted sum of rewards from different properties, and Section 3.3.2 shows the performance of MoleGuLAR in generating molecules with specified properties. To further validate the performance, the optimization was done from scratch across three runs with different random seeds, and the distribution of the properties of the generated molecules is discussed in the Supporting Information. Refer to Section S5 of the Supporting Information for the statistics and generated molecules from each of the tests in Section 3.3.

3.3.1. *Weighted Sum Reward.* Tests were done keeping the weights for each reward function equal to 1 in order to optimize the BA calculated using AutoDock along with target LogP = 2.5

generated, and their docking scores were calculated. After the completion of 175 iterations, 500 molecules were generated from the initial model and the optimized model. The distribution of the binding affinities given in Figures 5a and 5b shows a significant shift toward more desirable regions.

The above approach shows great performance in optimization for BA, but typical druglike molecules have constraints on other properties as well. Ideally, LogP should be between 0 and 5 for oral drugs, QED⁶⁰ close to 1, and TPSA < 90 \AA^2 , which are violated during single objective optimization (Figures 6a and

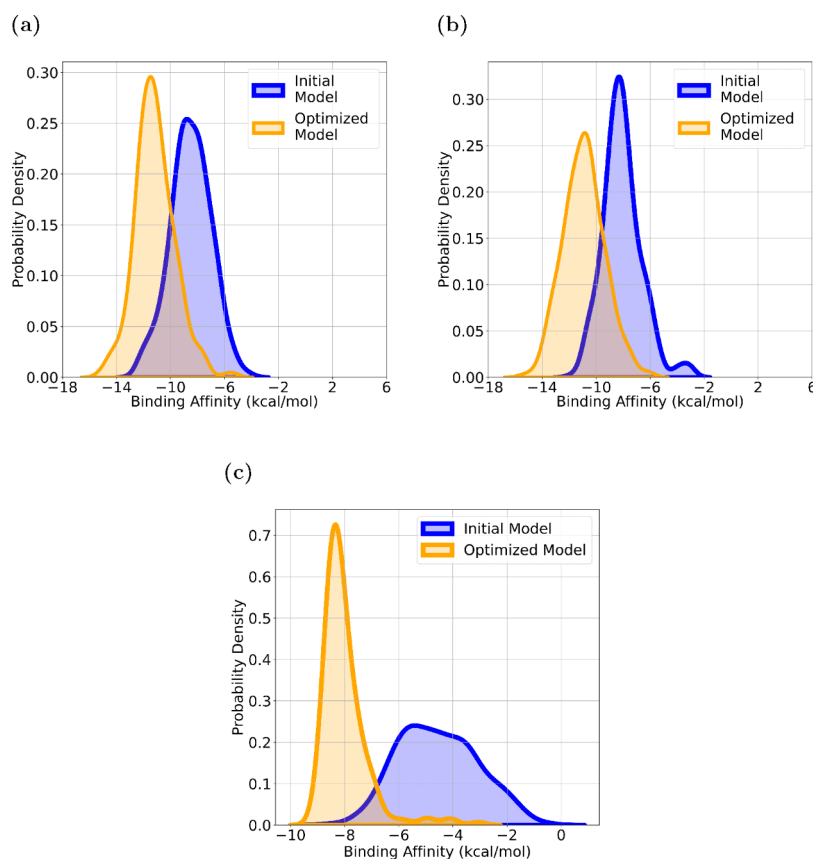


Figure 8. Distribution of BA of generated molecules before and after optimizing the generator for sum of rewards from target LogP = 2.5 and high BA calculated with (a) SARS-CoV-2 M_{pro} , (b) TTBK1, and (c) TTBK1 calculated using GIN.

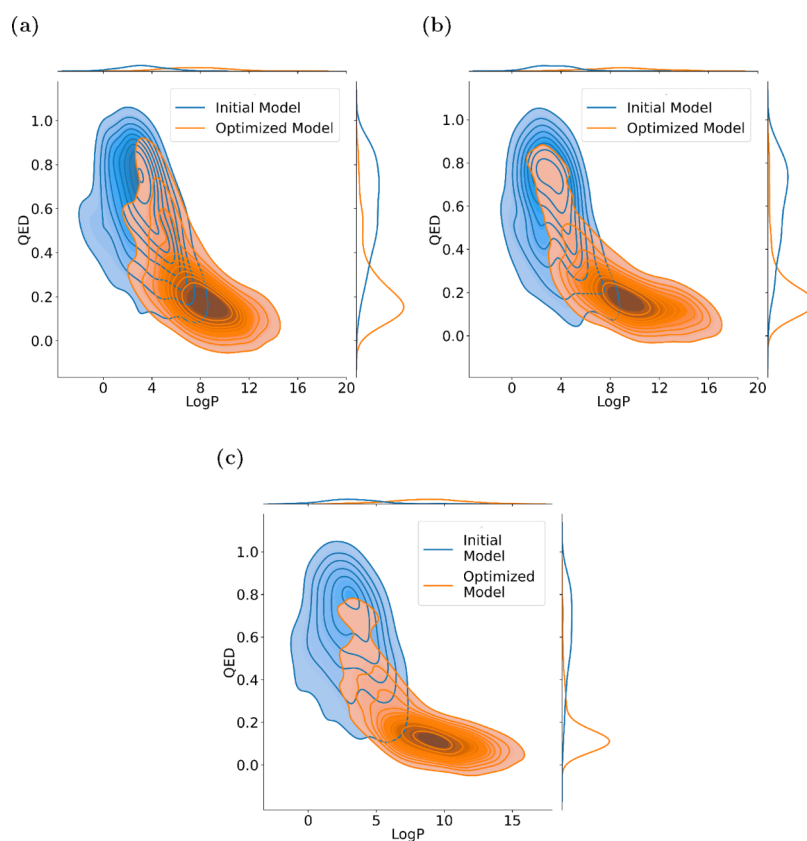


Figure 9. Joint distribution of LogP and QED of generated molecules before and after optimizing the generator for sum of rewards from target LogP = 2.5 and high BA with (a) SARS-CoV-2 M_{pro} , (b) TTBK1, and (c) TTBK1 calculated using GIN.

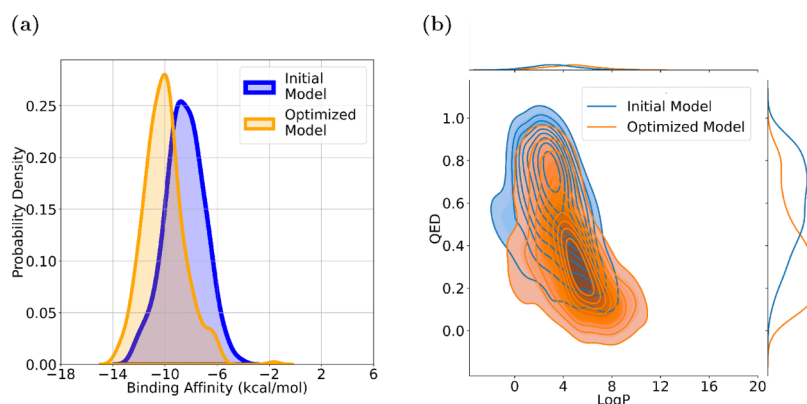


Figure 10. (a) Distribution of BA and (b) joint distribution of LogP and QED before and after optimizing the model for LogP = 2.5 and high BA with SARS-CoV-2 M_{pro} calculated using AutoDock by alternating rewards.

360 While there was improvement in the distribution of LogP in
 361 comparison to single objective optimization and BA in
 362 comparison to the initial model (Figure 8), the target was not
 363 achieved yet (Figure 9). A similar observation was seen when
 364 GIN was used for BA calculation instead of AutoDock in the
 365 same setting (Figures 8c and 9c). Further testing was done using
 366 the GIN BA predictor (Figure S16 of the Supporting
 367 Information), in which the generator was optimized to generate
 368 molecules with various simultaneous targets, i.e., LogP = 2.5,
 369 maximum QED, TPSA = 100 Å², and ΔG_{Hyd} = -10 kcal/mol,
 370 and the weighted sum of all rewards was taken. This worked well
 371 for BA, TPSA, and ΔG_{Hyd} but failed to achieve the target LogP
 372 and showed a very low QED. This is in agreement with the
 373 OptiMol pipeline by Boitreaud et al.⁴² who showed that

optimizing for BA led to a reduction in QED.⁴² In order to tackle
 this, we propose the following alternating rewards strategy for
 optimization.

3.3.2. *Alternating Rewards.* The pipeline's exceptional
 performance on single objective tasks helped formulate the
 strategy that only one objective be optimized at a time and the
 objective be changed at regular intervals. Taking the example of
 LogP and BA, initially the model moves to generating molecules
 with better BA, but after a few iterations, the reward is switched
 to optimize for LogP. Figure S19 of the Supporting Information
 shows the variation of BA with SARS-CoV-2 M_{pro} and LogP with
 the iterations. The generator is rewarded for BA during the
 iterations marked red and for LogP during the iterations marked
 blue, and it can be seen that when the reward switches the model

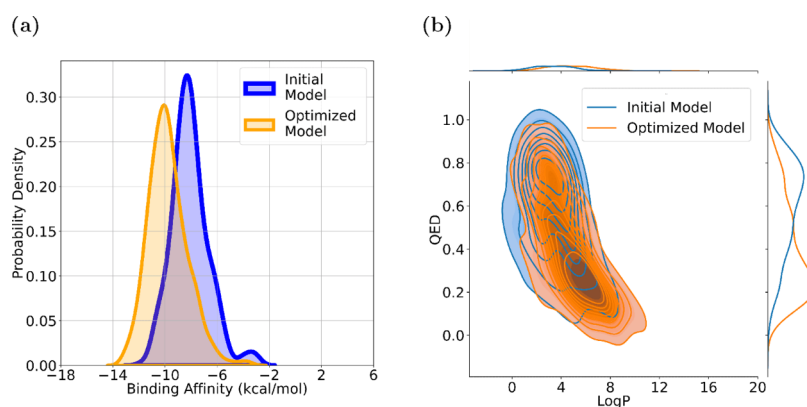


Figure 11. (a) Distribution of BA and (b) joint distribution of LogP and QED before and after optimizing the model for LogP = 2.5 and high BA with TTBK1 calculated using AutoDock by alternating rewards.

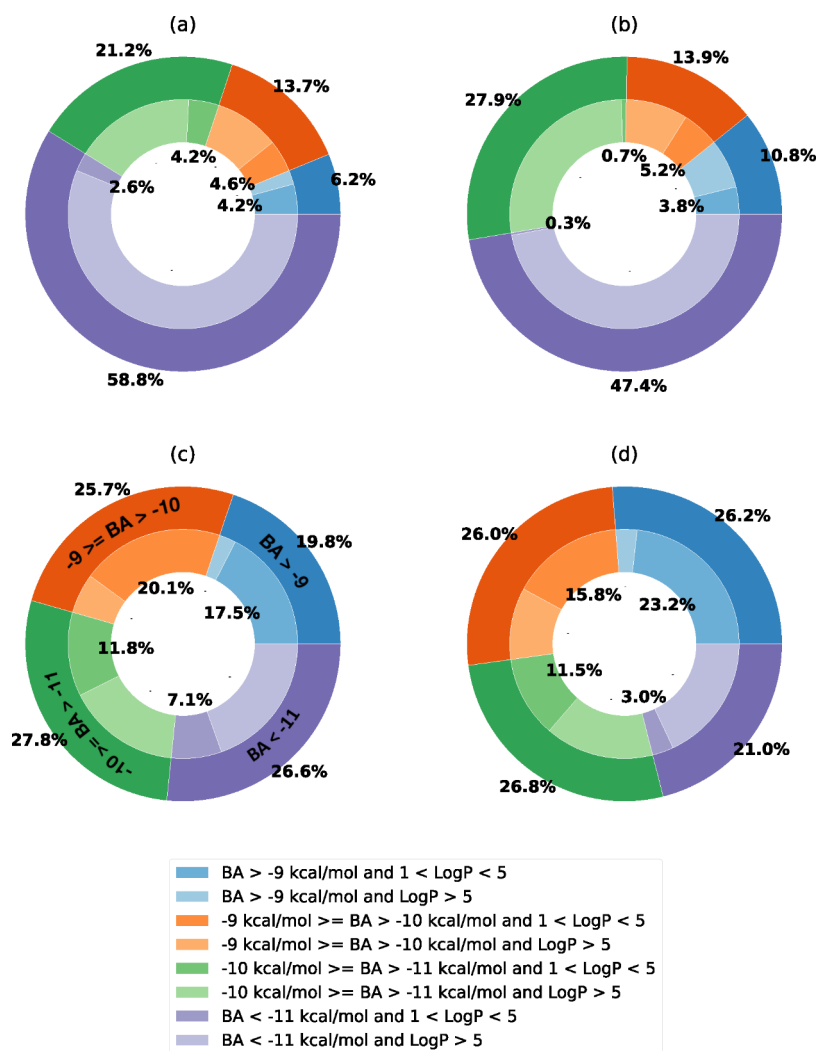


Figure 12. Fraction of molecules with BA and LogP in and out of desirable regions from those generated by the model optimized for (a) sum of rewards from BA with SARS-CoV-2 M_{pro} and LogP, (b) sum of rewards from BA with TTBK1 and LogP, (c) alternating rewards from BA with SARS-CoV-2 M_{pro} and LogP, and (d) alternating rewards from BA with TTBK1 and LogP. For legends of the outermost part of the pie chart, see (c), as they repeat in the rest of the pie charts.

388 is already sampling from the space with high BA and moves
 389 toward the region close to the target LogP and vice versa. Figures
 390 10 and 11 show the application of this strategy to SARS-CoV-2
 391 M_{pro} and TTBK1, respectively, using AutoDock. We can see a
 392 better distribution for BA as well as a significant overlap in the

most desirable and optimized regions of LogP and QED.
 Furthermore, Figure 12 shows the ability of the current strategy
 to consistently generate a higher percentage of hit molecules in
 comparison to the weighted sum approach.

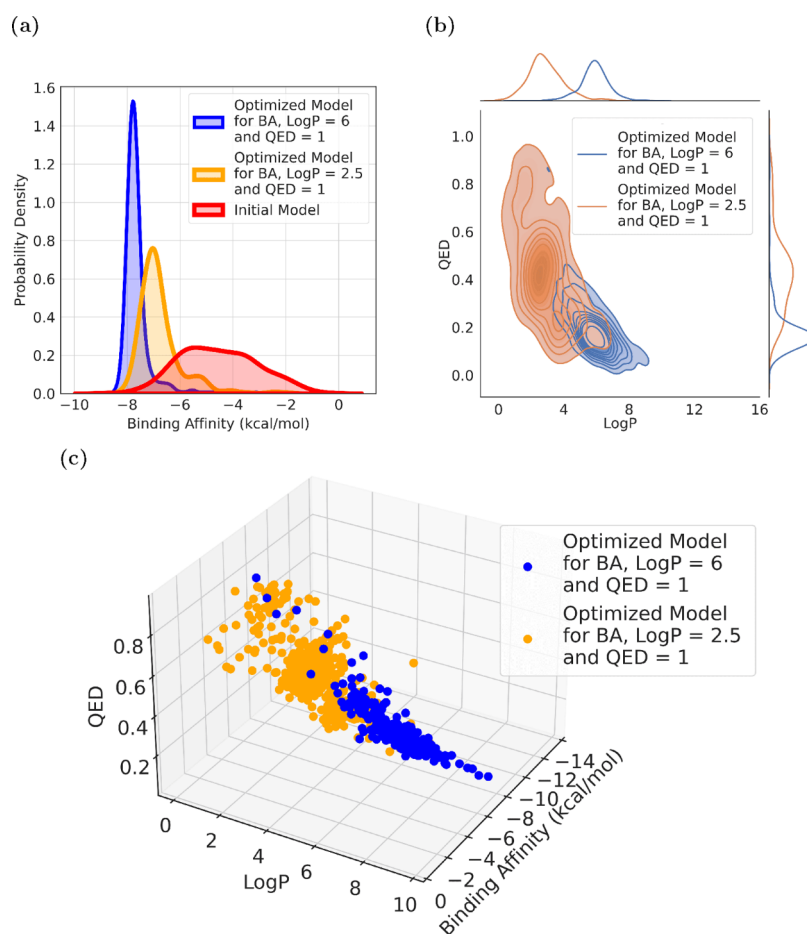


Figure 13. (a) Distribution of BA, (b) joint distribution of LogP and QED, and (c) 3D representation of properties of generated molecules before and after optimizing the generator for high BA with TTBK1 calculated using GIN.

Table 2. LogP and QED Targets along with Obtained Mean Values of BA, LogP, and QED of the Corresponding Generated Data as Well as the Best BA

target LogP	target QED	mean BA (kcal/mol)	best BA (kcal/mol)	mean LogP	mean QED
2.5	1	-6.76	-8.18	2.9	0.42
6	1	-7.64	-8.41	5.87	0.19

Table 3. TPSA and ΔG_{Hyd} Targets along with Mean Values of BA, TPSA, and ΔG_{Hyd} of the Corresponding Generated Data as Well as the Best BA

target TPSA (\AA^2)	target ΔG_{Hyd} (kcal/mol)	mean TPSA (\AA^2)	mean ΔG_{Hyd} (kcal/mol)	mean BA (kcal/mol)	best BA (kcal/mol)
70	-11	88.77	-10.13	-6.11	-8.36
120	-11	117.25	-10.75	-6.65	-8.32
70	-7	71.64	-7.42	-7.4	-8.90
120	-7	99.16	-8.42	-6.85	-8.64

different regions of the chemical space where molecules possess the desired properties.

A similar test was repeated with TPSA and ΔG_{Hyd} along with BA to see how the optimization strategy handles conflicting targets since higher the TPSA, the more negative the ΔG_{Hyd} . The target pairs are shown in Table 3, and the obtained results are discussed in Section S5.2.5 of the Supporting Information.

The best hits from 500 molecules generated from each of the aforementioned tests using alternating rewards are shown in Figure 14.

4. CONCLUSION

In this study, we propose MoleGuLAR, a pipeline for de novo generation of druglike molecules with high BA to novel targets along with other desirable properties using alternating rewards. Reinforcement learning is used to optimize the generator model weights to maximize the rewards obtained from calculated properties. A novel optimization strategy is also proposed for the multiobjective setup in which the reward function is switched to optimize for a different property at regular intervals instead of the conventional approach in which the sum of rewards from all properties is taken. We also show the performance of two ways of calculating BA, i.e., using AutoDock and using a predictor model, while also weighing the merits and demerits of both approaches as a part of the pipeline. Further work can include training the BA predictor models on the fly using techniques like active learning to make the pipeline more robust and efficient. The use of this architecture significantly reduces the number of

The GIN model was used for further testing. Different targets were kept for different properties to evaluate the model's capability of achieving all targets simultaneously. In Figure 13 and Table 2, it can be seen that the model is capable of generating molecules with high BA along with maximizing QED subject to the LogP being constrained to 2.5 and 6. Furthermore, in Figure 13c, there is a clear separation of the distributions in three dimensions showing the model's ability to navigate

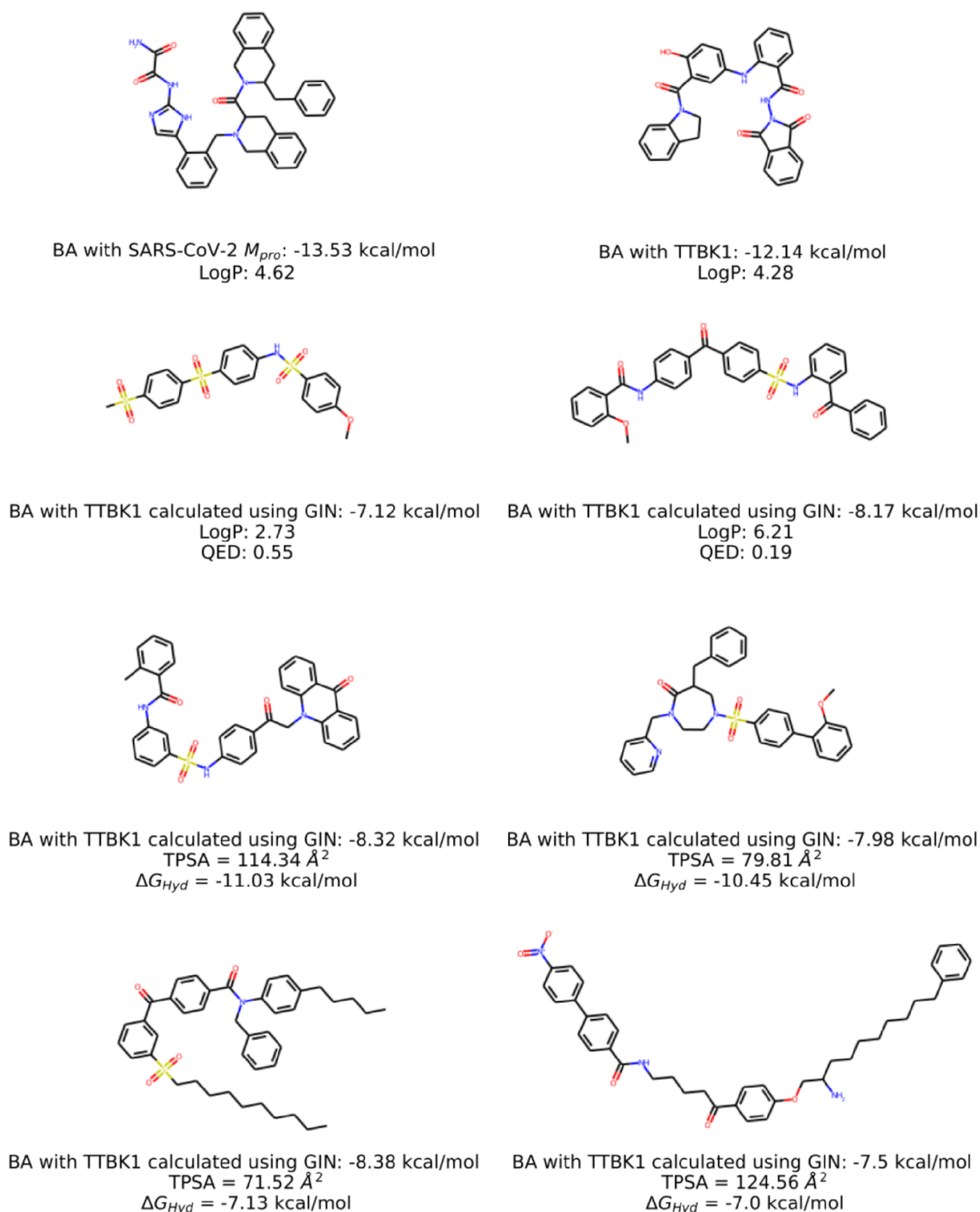


Figure 14. Top hits from 500 molecules generated after each test done using alternating rewards.

431 docking calculations required to identify potential drugs for a
 432 novel target removing a major bottleneck in the drug discovery
 433 process and can potentially be used to generate targeted drug
 434 libraries. We show that the alternating reward strategy is
 435 extremely robust in finding potential hits for the target across a
 436 wide set of target properties.

■ ASSOCIATED CONTENT

Supporting Information

The Supporting Information is available free of charge at
<https://pubs.acs.org/doi/10.1021/acs.jcim.1c01341>.

Performance of different machine-learning models for
 predicting BA, reward functions for each property, tables
 of mean values of respective properties of molecules

444 generated for each test along with structures of the top
445 hits, protein–ligand interactions, and supplementary
446 discussions and methods (PDF)

447 ■ AUTHOR INFORMATION

448 Corresponding Author

449 U. Deva Priyakumar – Center for Computational Natural
450 Sciences and Bioinformatics, International Institute of
451 Information Technology, Hyderabad 500 032, India;
452 orcid.org/0000-0001-7114-3955; Email: deva@iiit.ac.in

453 Authors

454 Manan Goel – Center for Computational Natural Sciences and
455 Bioinformatics, International Institute of Information
456 Technology, Hyderabad 500 032, India

457 Shampa Raghunathan – Center for Computational Natural
458 Sciences and Bioinformatics, International Institute of
459 Information Technology, Hyderabad 500 032, India; École
460 Centrale School of Engineering, Mahindra University,
461 Hyderabad 500 043, India

462 Siddhartha Laghuvarapu – Center for Computational Natural
463 Sciences and Bioinformatics, International Institute of
464 Information Technology, Hyderabad 500 032, India

465 Complete contact information is available at:

466 <https://pubs.acs.org/10.1021/acs.jcim.1c01341>

467 Notes

468 The authors declare no competing financial interest.

469 The source code, an extension of the work by Popova and co-
470 workers (<https://github.com/isayev/ReLeaSE>), along with
471 documentation and models optimized for each test are available
472 at <https://github.com/devalab/MoleGuLAR>.

473 ■ ACKNOWLEDGMENTS

474 We acknowledge the financial support through the DST-SERB
475 grant (no. CVD/2020/000 343), DST/WOS-A grant (no. SR/
476 WOS-A/CS-19/2 018 (G)), and IHUB-Data, IIIT Hyderabad.
477 This work was partially funded by Intel Corp. as part of its
478 Pandemic Response Technology Initiative (PRTI).

479 ■ REFERENCES

480 (1) Butler, K. T.; Davies, D. W.; Cartwright, H.; Isayev, O.; Walsh, A.
481 Machine learning for molecular and materials science. *Nature* **2018**,
482 *559*, 547–555.
483 (2) Sadowski, P.; Baldi, P. *Braverman Readings in Machine Learning*.
484 *Key Ideas from Inception to Current State*; Springer: New York, 2018; pp
485 269–297.
486 (3) Wu, Z.; Ramsundar, B.; Feinberg, E. N.; Gomes, J.; Geniesse, C.;
487 Pappu, A. S.; Leswing, K.; Pande, V. MoleculeNet: A benchmark for
488 molecular machine learning. *Chem. Sci.* **2018**, *9*, 513–530.
489 (4) Irwin, J. J.; Shoichet, B. K. ZINC—A free database of commercially
490 available compounds for virtual screening. *J. Chem. Inf. Model.* **2005**, *45*,
491 177–182.
492 (5) Mendez, D.; Gaulton, A.; Bento, A. P.; Chambers, J.; De Veij, M.;
493 Felix, E.; Magariños, M. P.; Mosquera, J. F.; Mutowo, P.; Nowotka, M.;
494 et al. ChEMBL: Towards direct deposition of bioassay data. *Nucleic*
495 *Acids Res.* **2019**, *47*, D930–D940.
496 (6) Senior, A. W.; Evans, R.; Jumper, J.; Kirkpatrick, J.; Sifre, L.; Green,
497 T.; Qin, C.; Židek, A.; Nelson, A. W. R.; Bridgland, A.; et al. Improved
498 protein structure prediction using potentials from deep learning. *Nature*
499 **2020**, *577*, 706–710.
500 (7) Segler, M. H.; Waller, M. P. Neural-symbolic machine learning for
501 retrosynthesis and reaction prediction. *Chem. - Eur. J.* **2017**, *23*, 5966–
502 5971.

(8) Schreck, J. S.; Coley, C. W.; Bishop, K. J. Learning retrosynthetic
503 planning through simulated experience. *ACS Cent. Sci.* **2019**, *5*, 970–
504 981.
505 (9) Gilmer, J.; Schoenholz, S. S.; Riley, P. F.; Vinyals, O.; Dahl, G. E.
506 *Neural message passing for quantum chemistry*; International Conference
507 on Machine Learning, 2017; pp 1263–1272.
508 (10) Szymański, P.; Markowicz, M.; Mikiciuk-Olasik, E. Adaptation of
509 high-throughput screening in drug discovery—toxicological screening
510 tests. *Int. J. Mol. Sci.* **2012**, *13*, 427–452.
511 (11) Broach, J. R.; Thorner, J. High-throughput screening for drug
512 discovery. *Nature* **1996**, *384*, 14–16.
513 (12) Maia, E. H. B.; Assis, L. C.; de Oliveira, T. A.; da Silva, A. M.;
514 Taranto, A. G. Structure-based virtual screening: From classical to
515 artificial intelligence. *Front. Chem.* **2020**, *8*.
516 (13) Sliwoski, G.; Kothiwale, S.; Meiler, J.; Lowe, E. W. Computa-
517 tional methods in drug discovery. *Pharmacol. Rev.* **2014**, *66*, 334–395.
518 (14) Aggarwal, R.; Gupta, A.; Chelur, V.; Jawahar, C. V.; Priyakumar,
519 U. D. DeepPocket: Ligand Binding Site Detection and Segmentation
520 using 3D Convolutional Neural Networks. *J. Chem. Inf. Model.* **2021**.
521 (15) Chelur, V.; Priyakumar, U. D. *BiRDS-Binding Residue Detection*
522 *from Protein Sequences using Deep ResNets*; **2021**.
523 (16) Mehta, S.; Laghuvarapu, S.; Pathak, Y.; Sethi, A.; Alvala, M.;
524 Priyakumar, U. D. MEMES: Machine learning framework for Enhanced
525 MolEcular Screening. *Chem. Sci.* **2021**, *12*, 11710.
526 (17) Lyu, J.; Wang, S.; Balias, T. E.; Singh, I.; Levit, A.; Moroz, Y. S.;
527 O'Meara, M. J.; Che, T.; Algae, E.; Tolmacheva, K.; Tolmachev, A. A.;
528 Shoichet, B. K.; Roth, B. L.; Irwin, J. J. Ultra-large library docking for
529 discovering new chemotypes. *Nature* **2019**, *566*, 224–229.
530 (18) Raymond, J.-L. The chemical space project. *Acc. Chem. Res.* **2015**,
531 *48*, 722–730.
532 (19) Gupta, A.; Müller, A. T.; Huisman, B. J.; Fuchs, J. A.; Schneider,
533 P.; Schneider, G. Generative recurrent networks for de novo drug
534 design. *Mol. Inf.* **2018**, *37*, 1700111.
535 (20) Weininger, D. SMILES, a chemical language and information
536 system. 1. Introduction to methodology and encoding rules. *J. Chem.*
537 *Inf. Model.* **1988**, *28*, 31–36.
538 (21) De Cao, N.; Kipf, T. MolGAN: An implicit generative model for
539 small molecular graphs; arXiv preprint; arXiv:1805.11973, **2018**.
540 (22) Prykhodko, O.; Johansson, S. V.; Kotsias, P.-C.; Arús-Pous, J.;
541 Bjerrum, E. J.; Engkvist, O.; Chen, H. A de novo molecular generation
542 method using latent vector based generative adversarial network. *J.*
543 *Cheminf.* **2019**, *11*, 1–13.
544 (23) Maziarika, E.; Pocha, A.; Kaczmarczyk, J.; Rataj, K.; Danel, T.;
545 Warchol, M. Mol-CycleGAN: a generative model for molecular
546 optimization. *J. Cheminf.* **2020**, *12*, 1–18.
547 (24) Jin, W.; Barzilay, R.; Jaakkola, T. *Junction tree variational*
548 *autoencoder for molecular graph generation*; International Conference on
549 Machine Learning, 2018; pp 2323–2332.
550 (25) Kusner, M. J.; Paige, B.; Hernández-Lobato, J. M. *Grammar*
551 *variational autoencoder*; International Conference on Machine Learn-
552 ing, 2017; pp 1945–1954.
553 (26) Griffiths, R.-R.; Hernández-Lobato, J. M. Constrained Bayesian
554 optimization for automatic chemical design using variational
555 autoencoders. *Chem. Sci.* **2020**, *11*, 577–586.
556 (27) Lim, J.; Ryu, S.; Kim, J. W.; Kim, W. Y. Molecular generative
557 model based on conditional variational autoencoder for de novo
558 molecular design. *J. Cheminf.* **2018**, *10*, 1–9.
559 (28) Kell, D. B.; Samanta, S.; Swainston, N. Deep learning and
560 generative methods in cheminformatics and chemical biology:
561 navigating small molecule space intelligently. *Biochem. J.* **2020**, *477*,
562 4559–4580.
563 (29) Jiménez-Luna, J.; Grisoni, F.; Weskamp, N.; Schneider, G.
564 Artificial intelligence in drug discovery: Recent advances and future
565 perspectives. *Expert Opin. Drug Discovery* **2021**, *16*, 949–959.
566 (30) Bagal, V.; Aggarwal, R.; Vinod, P.; Priyakumar, U. D. *LigGPT:*
567 *Molecular Generation using a Transformer-Decoder Model*; chemrxiv
568 preprint; **2021**. DOI: [10.26434/chemrxiv.14561901.v1](https://doi.org/10.26434/chemrxiv.14561901.v1).
569 (31) Winter, R.; Montanari, F.; Noé, F.; Clevert, D.-A. Learning
570 continuous and data-driven molecular descriptors by translating
571

- 572 equivalent chemical representations. *Chemical science* **2019**, *10*, 1692–
573 1701.
- 574 (32) Gómez-Bombarelli, R.; Wei, J. N.; Duvenaud, D.; Hernández-
575 Lobato, J. M.; Sánchez-Lengeling, B.; Sheberla, D.; Aguilera-
576 Iparraguirre, J.; Hirzel, T. D.; Adams, R. P.; Aspuru-Guzik, A.
577 Automatic chemical design using a data-driven continuous representa-
578 tion of molecules. *ACS Cent. Sci.* **2018**, *4*, 268–276.
- 579 (33) Koge, D.; Ono, N.; Huang, M.; Altaf-Ul-Amin, M.; Kanaya, S.
580 Embedding of Molecular Structure Using Molecular Hypergraph
581 Variational Autoencoder with Metric Learning. *Mol. Inf.* **2021**, *40*,
582 2000203.
- 583 (34) Tavakoli, M.; Baldi, P. Continuous Representation of Molecules
584 using Graph Variational Autoencoder. *Proceedings of the AAAI 2020*
585 *Spring Symposium on Combining Artificial Intelligence and Machine*
586 *Learning with Physical Sciences*, Stanford, CA, March 23–25, 2020.
- 587 (35) Winter, R.; Montanari, F.; Steffen, A.; Briem, H.; Noé, F.;
588 Clevert, D.-A. Efficient multi-objective molecular optimization in a
589 continuous latent space. *Chem. Sci.* **2019**, *10*, 8016–8024.
- 590 (36) Korovina, K.; Xu, S.; Kandasamy, K.; Neiswanger, W.; Poczos, B.;
591 Schneider, J.; Xing, E. *Chembo: Bayesian optimization of small organic*
592 *molecules with synthesizable recommendations*; International Conference
593 on Artificial Intelligence and Statistics, 2020; pp 3393–3403.
- 594 (37) Gao, K.; Nguyen, D. D.; Tu, M.; Wei, G.-W. Generative network
595 complex for the automated generation of drug-like molecules. *J. Chem.*
596 *Inf. Model.* **2020**, *60*, 5682–5698.
- 597 (38) Popova, M.; Isayev, O.; Tropsha, A. Deep reinforcement learning
598 for de novo drug design. *Sci. Adv.* **2018**, *4*, No. eaap7885.
- 599 (39) Guimaraes, G. L.; Sanchez-Lengeling, B.; Outeiral, C.; Farias, P.
600 L. C.; Aspuru-Guzik, A. *Objective-reinforced generative adversarial*
601 *networks (ORGAN) for sequence generation models*; arXiv preprint;
602 arXiv:1705.10843, **2017**.
- 603 (40) You, J.; Liu, B.; Ying, R.; Pande, V.; Leskovec, J. Graph
604 Convolutional Policy Network for Goal-Directed Molecular Graph
605 Generation. *Proceedings of the 32nd International Conference on Neural*
606 *Information Processing Systems*, Red Hook, NY, 2018; pp 6412–6422.
- 607 (41) Khemchandani, Y.; O'Hagan, S.; Samanta, S.; Swainston, N.;
608 Roberts, T. J.; Bollegala, D.; Kell, D. B. DeepGraphMolGen, a multi-
609 objective, computational strategy for generating molecules with
610 desirable properties: a graph convolution and reinforcement learning
611 approach. *J. Cheminf.* **2020**, *12*, 1–17.
- 612 (42) Boitreau, J.; Mallet, V.; Oliver, C.; Waldispohl, J. Optimol:
613 Optimization of binding affinities in chemical space for drug discovery.
614 *J. Chem. Inf. Model.* **2020**, *60*, 5658–5666.
- 615 (43) Born, J.; Manica, M.; Oskoei, A.; Cadow, J.; Markert, G.;
616 Martínez, M. R. PaccMannRL: De novo generation of hit-like
617 anticancer molecules from transcriptomic data via reinforcement
618 learning. *IScience* **2021**, *24*, 102269–269.
- 619 (44) Born, J.; Manica, M.; Cadow, J.; Markert, G.; Mill, N. A.;
620 Filipavicius, M.; Janakarajan, N.; Cardinale, A.; Laino, T.; Martínez, M.
621 R. Data-driven molecular design for discovery and synthesis of novel
622 ligands: a case study on SARS-CoV-2. *Mach. Learn.: Sci. Technol.* **2021**,
623 *2*, 025–024.
- 624 (45) Bung, N.; Krishnan, S. R.; Bulusu, G.; Roy, A. De novo design of
625 new chemical entities for SARS-CoV-2 using artificial intelligence.
626 *Future Med. Chem.* **2021**, *13*, 575–585.
- 627 (46) Krishnan, S. R.; Bung, N.; Bulusu, G.; Roy, A. Accelerating de
628 novo drug design against novel proteins using deep learning. *J. Chem.*
629 *Inf. Model.* **2021**, *61*, 621–630.
- 630 (47) Joulin, A.; Mikolov, T. Inferring Algorithmic Patterns with Stack-
631 Augmented Recurrent Nets. *Proceedings of the 28th International*
632 *Conference on Neural Information Processing Systems*, Cambridge, MA,
633 2015; Vol. 1, pp 190–198.
- 634 (48) Landrum, G. *RDKit: A software suite for cheminformatics,*
635 *computational chemistry, and predictive modeling*; 2013.
- 636 (49) Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T.
637 N.; Weissig, H.; Shindyalov, I. N.; Bourne, P. E. The protein data bank.
638 *Nucleic acids research* **2000**, *28*, 235–242.
- 639 (50) Morris, G. M.; Huey, R.; Lindstrom, W.; Sanner, M. F.; Belew, R.
640 K.; Goodsell, D. S.; Olson, A. J. AutoDock4 and AutoDockTools4:
Automated docking with selective receptor flexibility. *J. Comput. Chem.* **2009**, *30*, 2785–2791.
- (51) Santos-Martins, D.; Solis-Vasquez, L.; Tillack, A. F.; Sanner, M.
643 F.; Koch, A.; Forli, S. Accelerating AutoDock4 with GPUs and gradient-
644 based local search. *J. Chem. Theory Comput.* **2021**, *17*, 1060–1073.
- (52) HTS Collection. [https://enamine.net/compound-collections/
645 screening-collection/hts-collection](https://enamine.net/compound-collections/screening-collection/hts-collection) (accessed on 11/23/2021).
- (53) Jaeger, S.; Fulle, S.; Turk, S. Mol2vec: unsupervised machine
648 learning approach with chemical intuition. *J. Chem. Inf. Model.* **2018**, *58*,
649 27–35.
- (54) Xu, K.; Hu, W.; Leskovec, J.; Jegelka, S. *How Powerful are Graph*
651 *Neural Networks?* 7th International Conference on Learning
652 Representations, ICLR 2019, New Orleans, LA, May 6–9, 2019.
- (55) Hu, W.; Liu, B.; Gomes, J.; Zitnik, M.; Liang, P.; Pande, V. S.;
654 Leskovec, J. *Strategies for Pre-training Graph Neural Networks*; 8th
655 International Conference on Learning Representations, ICLR 2020,
656 Addis Ababa, Ethiopia, April 26–30, 2020.
- (56) Pathak, Y.; Mehta, S.; Priyakumar, U. D. Learning Atomic
658 Interactions through Solvation Free Energy Prediction Using Graph
659 Neural Networks. *J. Chem. Inf. Model.* **2021**, *61*, 689–698.
- (57) Williams, R. J. Simple statistical gradient-following algorithms for
661 connectionist reinforcement learning. *Machine Learning* **1992**, *8*, 229–
662 256.
- (58) Zhang, L.; Lin, D.; Sun, X.; Curth, U.; Drosten, C.; Sauerhering,
664 L.; Becker, S.; Rox, K.; Hilgenfeld, R. Crystal structure of SARS-CoV-2
665 main protease provides a basis for design of improved α -ketoamide
666 inhibitors. *Science* **2020**, *368*, 409–412.
- (59) Sato, S.; Cerny, R. L.; Buescher, J. L.; Ikezu, T. Tau-tubulin
668 kinase 1 (TTBK1), a neuron-specific tau kinase candidate, is involved in
669 tau phosphorylation and aggregation. *J. Neurochem.* **2006**, *98*, 1573–
670 1584.
- (60) Bickerton, G. R.; Paolini, G. V.; Besnard, J.; Muresan, S.;
672 Hopkins, A. L. Quantifying the chemical beauty of drugs. *Nat. Chem.* **2012**, *4*, 90–98.