

# AN EVALUATION FRAMEWORK FOR THE OBJECTIVE FUNCTIONS OF *de novo* DRUG DESIGN BENCHMARKS

**Austin Tripp**

University of Cambridge  
ajt212@cam.ac.uk

**Wenlin Chen**

University of Cambridge  
MPI for Intelligent Systems  
wc337@cam.ac.uk

**José Miguel Hernández-Lobato**

University of Cambridge  
jmh233@cam.ac.uk

## ABSTRACT

*De novo* drug design has recently received increasing attention from the machine learning community. It is important that the field is aware of the actual goals and challenges of drug design and the roles that *de novo* molecule design algorithms could play in accelerating the process, so that algorithms can be evaluated in a way that reflects how they would be applied in real drug design scenarios. In this paper, we propose a framework for critically assessing the merits of benchmarks, and argue that most of the existing *de novo* drug design benchmark functions are either highly unrealistic or depend upon a surrogate model whose performance is not well characterized. In order for the field to achieve its long-term goals, we recommend that poor benchmarks (especially logP and QED) be deprecated in favour of better benchmarks. We hope that our proposed framework can play a part in developing new *de novo* drug design benchmarks that are more realistic and ideally incorporate the intrinsic goals of drug design.

## 1 INTRODUCTION

*De novo* molecular design, also called molecular optimization, is the problem of producing molecules with desirable properties. Mathematically it is most often formulated as maximizing an objective function which encodes the desirability of a molecule. A growing number of papers at major machine learning conferences address this problem by proposing general algorithms to optimize functions over molecule space (see Appendix A and Elton et al. (2019)).

Unfortunately, almost none of these works actually evaluate their algorithms on tasks that could reasonably be considered real-world drug-design problems. This is understandable: nearly all molecular properties of practical interest are quantities derived from wet-lab experiments or computationally-intensive simulations, making it very impractical (if not impossible) for machine learning researchers to evaluate a real-world objective function on a novel molecule. It is presumably for this reason that machine learning papers optimize cheap, computer-based objective functions to demonstrate the performance of their algorithms, such as penalized logP (Kusner et al., 2017; Gómez-Bombarelli et al., 2018), QED (Bickerton et al., 2012), or the GuacaMol benchmark suite (Brown et al., 2019). Over time, these benchmark functions have become somewhat standardized, and researchers submitting work to top-tier machine learning conferences are expected to test their algorithms on at least some of these benchmarks.

Amidst all these standardized benchmarks it is easy to lose sight of this field’s primary long-term goal: producing algorithms for real-world drug design problems. We think the time is right for the community to step back and ask itself: “will getting higher scores on these benchmarks *really* lead to meaningful advancements in the field, or just produce a lot of conference papers?” We attempt to partially answer this question for *de novo* molecular design benchmarks by focusing on the relevance of their objective functions.

There are different types of objective functions whose usefulness stems from different sources. In Table 1 we propose a simple framework to classify objective functions into three different types: intrinsic objectives, approximations of intrinsic objectives, and proxy objectives (i.e., everything else). Benchmarks maximizing intrinsic objectives should be uncontroversially useful because they

TYPE	DESCRIPTION	EXAMPLE	HOW TO EVALUATE
<b>Intrinsic Objective</b>	An objective whose optimization is intrinsically useful.	Experimental binding to protein	Always useful
<b>Approximate Intrinsic Objective</b>	A direct approximation or simplification of an intrinsic objective.	Simulated binding energy	Quality of approximation
<b>Proxy Objective</b>	An objective whose optimization does not resemble or approximate an intrinsic goal but is optimized anyway.	Quantitative estimate of drug-likeness (QED)	Resemblance to some intrinsic objective

Table 1: Our proposed taxonomy of objective functions.

represent a real-world goal; however the high cost of real-world objective functions makes such benchmarks rare in practice. Benchmarks maximizing an *approximation* to an intrinsic objective could still be useful, but this utility will naturally diminish as the accuracy of the approximation decreases. Therefore the approximation quality is the most natural criterion for evaluating the utility of such benchmarks. However, the case *proxy* objectives is less clear: if an objective is not even an approximation of something that researchers value, why should researchers care about maximizing it? Given that researchers in machine learning mainly use benchmarks to determine which algorithms and techniques are worth further study, we believe that the merit of a proxy objective stems from its resemblance to an intrinsic objective. If this resemblance is very strong then there is good reason to suspect that strong performance on a proxy objective will translate to strong performance on a real-world objective, and it is therefore sensible to use the proxy objective as a benchmark. Conversely, if a proxy objective does not resemble any intrinsic objective then we believe using it as a benchmark provides no value to the research community.

The main contribution of this paper is to use the framework in Table 1 to evaluate some commonly-used benchmarks for machine learning in drug design. We start by giving an overview of some intrinsic goals and challenges of drug design in section 2, which provides important context for the benchmarks we discuss. We then apply our evaluation framework to some common *de novo* drug design benchmarks in section 3, and explain why many popular benchmarks are less useful than many researchers may expect. We end in section 4 with some thoughts about what actions researchers can take in the short and medium term to ensure that the field stays focused on its long-term goals. We hope that this paper will serve as a useful reference and a call to action for the machine learning in drug design community.

## 2 WHAT ARE THE “REAL” GOALS OF DRUG DESIGN AND HOW COULD *de novo* DESIGN ALGORITHMS HELP?

This section attempts to help the reader understand what the intrinsic goals actually are for drug discovery by providing a high level overview of drug design. There are many substances that can be used as medicines, but the field of “drug design” generally refers more specifically to finding small organic molecules which interact with specific biomolecules in the body (e.g. proteins). Typically the drug design process involves 1) identifying a biomolecule of interest, 2) finding a small molecule which binds to this biomolecule, and 3) performing larger tests in petri dishes, animals, or humans to determine whether the small molecule actually achieves the desired therapeutic effect.<sup>1</sup> Although innovations which contribute to any of these steps would be useful for drug design, algorithms for *de novo* molecular design would presumably be applicable to step 2 of this process because it is the only step that actually involves designing molecules.

<sup>1</sup>There are a variety of reasons why a molecule might not have the desired therapeutic effect. For example, the drug may degrade in the body, be excreted by the body too quickly, or have side effects.

If this problem is framed as optimization, there are many possible choices for the objective function, such as the empirical binding strength (to be maximized), the concentration required for the inhibition of protein activity (to be minimized), or a computational estimate of the binding free energy (to be minimized). It may also be sensible to incorporate constraints into the objective, either directly related to the protein interaction or related to other important properties (such as Lipinski’s rules for the drug to be absorbed orally (Lipinski et al., 1997)). Although every objective function is different, decades of drug discovery research have demonstrated that practical objective functions (to be maximized) often have the following qualities:

1. **Infrequent optima.** The overwhelming majority of molecules have low (i.e., poor) objective function values. “Random sampling” techniques have a very low rate of success, as evidenced by previous results in high-throughput screening (Bender et al., 2008).
2. **Multiple modes.** There are often several structurally distinct clusters of molecules with high objective function values. For example, Cephalosporins and Penicillins are two antibiotic drug classes which have only a small substructure in common, but both affect the same target protein.
3. **Complex relationship between molecular structure and objective function value.** For example, having a particular set of substructures may be a necessary but not sufficient condition for a molecule to have a high objective function value.<sup>2</sup>
4. **Variable “smoothness”.** Although objective functions are often “smooth” (i.e., small changes in molecular structure result in a small change in objective function values), this is usually not the case globally: there are some regions of molecular space where small changes in structure cause large changes in objective function values. This is often called *activity cliffs* in the drug design community (Maggiora, 2006).
5. **Noisy and expensive evaluations.** Each function evaluation is typically quite costly and has associated measurement noise. This noise can be reduced by repeated measurements or using more precise instruments, but at a much higher cost. The noise may also be non-uniform and have a non-zero mean.

Unfortunately, all of the above properties make “real-world” drug discovery objectives difficult to optimize. This is likely why drug design is considered a hard problem and why many classical optimization techniques have not performed well in practice: if the problem were easier there would likely not be much interest in the creation of novel algorithms and methodologies. We believe that the primary challenge for *de novo* design algorithms is to work well *in spite* of these difficulties.

### 3 COMMENTS ON SPECIFIC *de novo* DRUG DESIGN BENCHMARKS

In this section, we apply our framework to judge the merits of some commonly-used drug discovery benchmarks in the machine learning literature.

#### 3.1 PENALIZED LOGP MAXIMIZATION

The octanol-water partition coefficient (sometimes called  $P$ ) is the ratio between a molecule’s solubility in alcohol and its solubility in water. This value is of great interest to drug design because it influences how a molecule will be transported and absorbed by the body (which has both water and oily tissues). The most common real-world usage of logP is as preliminary filter to decide whether a drug candidate should be discarded: Lipinski et al. (1997) suggested in a highly cited paper that molecules with  $\log P > 5$  are unlikely to be absorbed by the body and are therefore not worth testing. Although not absolute, chemists generally accept this as a rule of thumb.

Many machine learning papers utilize benchmarks which involve a computational approximation of logP (Wildman & Crippen, 1999). The most common benchmarks in this class involve maximizing logP subject to a penalty, usually synthetic accessibility score (SAS) (Ertl & Schuffenhauer, 2009; Kusner et al., 2017; Gómez-Bombarelli et al., 2018) or similarity to a held-out target molecule (Jin

<sup>2</sup>For example, the 3D geometry of the molecule might need to be such that the substructures have a particular orientation. Alternatively, there may be other substructures which block the active substructures, so a molecule must *not* have some substructures in addition to having others.

et al., 2018). In either case, the objective rewards algorithms which produce molecules whose logP is much greater than 5, and are therefore very unlikely to be useful drug candidates.

Clearly, producing molecules with high logP values is not an intrinsic goal in drug design, nor does it directly approximate an intrinsic goal. Therefore under our framework penalized logP is a proxy objective and should be evaluated based on its resemblance to an intrinsic objective, which we believe to be negligible. The functional form of the logP approximation from Wildman & Crippen (1999) is a weighted sum over the atoms in a molecule, with the weights coming from a table with 68 mutually exclusive atom types. This means that the logP value can always be increased by adding more atoms with positive logP weights, such as carbon and chlorine, with no complex interactions between the different atoms. Consequently, there are no local optima (the logP of every molecule can be increased by adding extra atoms), and almost every large molecule will have a high logP value. This makes optimizing logP very easy.

These attributes are not fundamentally changed by the various penalties present in popular objectives: the synthetic accessibility penalty has a maximum value of 10 (and therefore just decreases the logP score of large molecules by a constant value), while the similarity penalty is a binary constraint (and therefore does not modify the objective so long as the constraints are satisfied). Overall, because penalized logP maximization is not an intrinsically valuable objective and possesses none of the difficult features of real-world objective functions described in section 2, we believe that it is a poor choice of benchmark and provides little or no value to the community.

### 3.2 QED

The quantitative estimate of drug-likeness (QED) is a heuristic metric which, as the name suggests, tries to quantify how drug-like a molecule is (Bickerton et al., 2012). Unlike the name might suggest however, a high QED score does not imply that a molecule is a promising drug candidate: it merely shows that it lacks a number of “red flags” that chemists generally think are not promising. The converse is also not true: many real approved drug molecules have low QED scores. Because of this, in practice QED is used as a metric for screening rather than an explicit optimization objective as is seen in machine learning papers. We therefore believe that QED should be classified as a “proxy objective” in our framework.

Similar to that of penalized logP, the relationship between molecular structure and QED value is fairly simple, suggesting that it might be very easy to optimize. In a recent workshop paper Tripp et al. (2021) found that a large number of algorithms (including random search) are able to produce molecules with a QED of 0.948 after a fairly small number of iterations, and that no paper to date has reported a molecule with a QED  $> 0.948$  (implying that this is likely the global maximum). It is therefore reasonable to conclude that QED maximization is a poor proxy benchmark, mainly because it is too simple to meaningfully distinguish different algorithms.

### 3.3 GUACAMOL

The GuacaMol benchmark suite contains 20 objective functions based on fingerprint similarity (Brown et al., 2019). To maximize each objective an algorithm must produce a molecule with high similarity to an existing drug molecule, but with some slightly altered properties (e.g. a different number of rings). The objectives have some of the properties discussed in section 2 and were designed to resemble a class of objective functions commonly seen in real-world drug design: producing modified versions of an initial “lead” molecule. Therefore, according to our framework these objectives should be considered as high-quality proxy objectives.

### 3.4 PREDICTED JNK3 AND GNK3 $\beta$ INHIBITION

Recently a number of papers perform experiments maximizing the predicted activity of a molecule against the targets JNK3 and GNK3 $\beta$ , with the predictor typically being a pre-trained random forest model (Li et al., 2018; Jin et al., 2020b). JNK3 and GNK3 $\beta$  are real proteins with pharmaceutical relevance and could therefore be considered intrinsic objectives. This makes optimizing the predictions of a random forest model an approximate intrinsic objective under our proposed framework. Accordingly, the utility of this objective function can be judged by the accuracy of the random forest predictor. Unfortunately, to our knowledge there have been no thorough investigations of these

particular random forest models besides the initial report of an AUROC score in (Jin et al., 2020b) and an examination of the compounds produced. Given that the random forest models were trained on less than  $10^4$  data points, it is plausible that the accuracy could be very poor in regions of chemical space far away from the training data. Ultimately this makes the usefulness of this objective unclear.

### 3.5 DOCKSTRING

The DOCKSTRING package (García-Ortegón et al., 2021) provides several interesting functions based on docking scores, which result from an explicit simulation of molecule-protein interactions using AutoDock Vina (Trott & Olson, 2010). The objectives are all motivated by real-world drug discovery problems, but are defined using AutoDock Vina’s binding scores, and are therefore approximate intrinsic objectives under our framework.

The merit of these benchmarks depends on the quality of AutoDock Vina, of which reports are mixed. On one hand, García-Ortegón et al. (2021) report only a modest (but still non-trivial) correlation between docking scores and experimentally measured binding affinities, suggesting that the approximation quality may not be very good. On the other hand, the DOCKSTRING objectives use a QED penalty to correct for some known biases in AutoDock Vina’s scoring algorithms, which has the potential to partially mitigate the low approximation accuracy. Overall we believe that these benchmarks are still sufficiently challenging as to be useful for the community, but acknowledge that this may be disproved in the future. It must also be noted that the DOCKSTRING objectives are much more computationally intensive than all other objective functions described so far.

## 4 CONCLUSION: WHERE SHOULD WE GO FROM HERE?

In this paper, we gave an overview of the goals of drug design, presented a framework to evaluate the merits of benchmarks in relation to these high level goals, then applied the framework to some common *de novo* design objectives. Unsurprisingly, none of the benchmarks studied used an intrinsic goal as an objective function. Two benchmarks, DOCKSTRING and JNK3/GNK3 $\beta$  directly approximate intrinsic objectives, although the global accuracy of these approximations is unknown, making it unclear how useful these objectives are in practice. The remaining three benchmarks could be considered “proxy objectives”. Of these, we found that penalized logP and QED are very simple to optimize and thereby do not resemble realistic benchmark functions, implying that their value to the community is limited. The GuacaMol objectives on the other hand do closely resemble realistic problems in drug design and therefore are likely useful benchmarks.

Overall, it seems that most of the benchmarks studied are either not useful or have unclear utility: a finding that we find unsatisfying. Given this, we have several suggestions for the community to ensure that the field stays focused on its long-term goals. In the short term, we believe researchers should be more discerning about the benchmarks that they use in papers. Firstly, it would be helpful if benchmarks such as penalized logP and QED maximization were avoided by the community, as they are very easy and their further use encourages even more papers to use them. The GuacaMol benchmarks are an excellent substitute with similar computational requirements. Secondly, while creating new benchmarks can be valuable and important, it is vital that the relationship between the proposed benchmark and real-world problems be explicitly considered. Far too many papers create arbitrary-seeming functions without justifying their utility, seemingly unaware that the no free lunch theorem implies that optimizing arbitrary functions is uninformative. Here, our framework may prove helpful (although it is far from the only way to view objective functions).

In the longer term, we believe that the community would benefit from a concerted effort to discuss the field’s medium and long term goals, then propose benchmarks to very explicitly measure progress towards these goals. Firstly, it would be extremely valuable to further investigate the accuracy of the approximations in benchmarks like JNK3/GNK3 $\beta$  and DOCKSTRING, to better understand the shortcomings of these existing benchmarks. Secondly, there is still a lot of value in proposing new objective functions to optimize, especially for different sub-problems in drug design. For example, although the task of “lead optimization” is approximated fairly well by the GuacaMol benchmark, the task of finding novel promising lead molecules (i.e. active molecules substantially different from known active molecules) is relatively different and arguably much more challenging. It would be valuable to create objective functions which test this. Finally, we think it is important to also consider

new paradigms for evaluating *de novo* design algorithms. While maximizing an objective function is a fairly general and flexible problem formulation, it fails to capture many aspects present in real world problems such as multiple sources of information, multiple objectives, and heterogeneous experimental noise. Many practical problems could be better formulated in other ways, and this should be investigated by the community.

Overall we are optimistic about the prospects of machine learning and *de novo* molecular design, but believe that the community must be vigilant to ensure that its success metrics are aligned towards achieving long-term goals, even if these goals are not currently measurable or achievable.

## ACKNOWLEDGEMENTS

We thank an anonymous reviewer for a thoughtful critique of this manuscript. We will incorporate this feedback into the next version of this text.

## REFERENCES

- Sungsoo Ahn, Junsu Kim, Hankook Lee, and Jinwoo Shin. Guiding deep molecular optimization with genetic exploration. *Advances in neural information processing systems*, 33:12008–12021, 2020.
- Sungsoo Ahn, Binghong Chen, Tianzhe Wang, and Le Song. Spanning tree-based graph generation for molecules. In *International Conference on Learning Representations*, 2022. URL [https://openreview.net/forum?id=w60btE\\_8T2m](https://openreview.net/forum?id=w60btE_8T2m).
- Andreas Bender, Dejan Bojanic, John W Davies, Thomas J Crisman, Dmitri Mikhailov, Josef Scheiber, Jeremy L Jenkins, Zhan Deng, W Adam G Hill, Maxim Popov, et al. Which aspects of hits are empirically correlated with downstream success? *Current Opinion in Drug Discovery and Development*, 11(3):327, 2008.
- Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, and Yoshua Bengio. Flow network based generative models for non-iterative diverse candidate generation. *Advances in Neural Information Processing Systems*, 34, 2021.
- G Richard Bickerton, Gaia V Paolini, Jérémy Besnard, Sorel Muresan, and Andrew L Hopkins. Quantifying the chemical beauty of drugs. *Nat. Chem.*, 4(2):90–98, 2012.
- John Bradshaw, Brooks Paige, Matt J Kusner, Marwin Segler, and José Miguel Hernández-Lobato. A model to search for synthesizable molecules. *Advances in Neural Information Processing Systems*, 32, 2019.
- John Bradshaw, Brooks Paige, Matt J Kusner, Marwin Segler, and José Miguel Hernández-Lobato. Barking up the right tree: an approach to search over molecule synthesis dags. *Advances in Neural Information Processing Systems*, 33:6852–6866, 2020.
- Nathan Brown, Marco Fiscato, Marwin HS Segler, and Alain C Vaucher. Guacamol: benchmarking models for de novo molecular design. *J. Chem. Inf. Model.*, 59(3):1096–1108, 2019.
- Vijil Chenthamarakshan, Payel Das, Samuel Hoffman, Hendrik Strobelt, Inkit Padhi, Kar Wai Lim, Benjamin Hoover, Matteo Manica, Jannis Born, Teodoro Laino, and Aleksandra Majsilovic. Cogmol: Target-specific and selective drug design for covid-19 using deep generative models. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 4320–4332. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/2d16ad1968844a4300e9a490588ff9f8-Paper.pdf>.
- Daniel C Elton, Zois Boukouvalas, Mark D Fuge, and Peter W Chung. Deep learning for molecular design—a review of the state of the art. *Molecular Systems Design & Engineering*, 4(4):828–849, 2019.

- Peter Ertl and Ansgar Schuffenhauer. Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions. *Journal of cheminformatics*, 1(1):1–11, 2009.
- Tianfan Fu, Wenhao Gao, Cao Xiao, Jacob Yasonik, Connor W. Coley, and Jimeng Sun. Differentiable scaffolding tree for molecule optimization. In *International Conference on Learning Representations*, 2022. URL [https://openreview.net/forum?id=w\\_drCosT76](https://openreview.net/forum?id=w_drCosT76).
- Wenhao Gao, Rocío Mercado, and Connor W. Coley. Amortized tree generation for bottom-up synthesis planning and synthesizable molecular design. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=FRxhHdnxt1>.
- Miguel García-Ortegón, Gregor NC Simm, Austin J Tripp, José Miguel Hernández-Lobato, Andreas Bender, and Sergio Bacallado. Dockstring: easy molecular docking yields better benchmarks for ligand design. *arXiv preprint arXiv:2110.15486*, 2021.
- Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato, Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. *ACS Cent. Sci.*, 4(2):268–276, 2018.
- Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Junction tree variational autoencoder for molecular graph generation. In *International conference on machine learning*, pp. 2323–2332. PMLR, 2018.
- Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Hierarchical generation of molecular graphs using structural motifs. In *International Conference on Machine Learning*, pp. 4839–4848. PMLR, 2020a.
- Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Multi-objective molecule generation using interpretable substructures. In *International conference on machine learning*, pp. 4849–4859. PMLR, 2020b.
- Matt J Kusner, Brooks Paige, and José Miguel Hernández-Lobato. Grammar variational autoencoder. In *International Conference on Machine Learning*, pp. 1945–1954. PMLR, 2017.
- Yibo Li, Liangren Zhang, and Zhenming Liu. Multi-objective de novo drug design with conditional graph generative model. *Journal of cheminformatics*, 10(1):1–24, 2018.
- Christopher A Lipinski, Franco Lombardo, Beryl W Dominy, and Paul J Feeney. Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. *Advanced drug delivery reviews*, 23(1-3):3–25, 1997.
- Qi Liu, Miltiadis Allamanis, Marc Brockschmidt, and Alexander Gaunt. Constrained graph variational autoencoders for molecule design. *Advances in Neural Information Processing Systems*, 31:7795–7804, 2018.
- Xianggen Liu, Qiang Liu, Sen Song, and Jian Peng. A chance-constrained generative framework for sequence optimization. In *International Conference on Machine Learning*, pp. 6271–6281. PMLR, 2020.
- Youzhi Luo, Keqiang Yan, and Shuiwang Ji. Graphdf: A discrete flow model for molecular graph generation. In *International Conference on Machine Learning*, pp. 7192–7203. PMLR, 2021.
- Gerald M Maggiora. On outliers and activity cliffs why qsar often disappoints. *Journal of chemical information and modeling*, 46(4):1535–1535, 2006.
- Krzysztof Maziarz, Henry Richard Jackson-Flux, Pashmina Cameron, Finton Sirockin, Nadine Schneider, Nikolaus Stiefl, Marwin Segler, and Marc Brockschmidt. Learning to extend molecular scaffolds with structural motifs. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=ZTsoE8G3GG>.

- Amina Mollaysa, Brooks Paige, and Alexandros Kalousis. Goal-directed generation of discrete structures with conditional generative models. *Advances in Neural Information Processing Systems*, 33:21923–21933, 2020.
- Henry Moss, David Leslie, Daniel Beck, Javier Gonzalez, and Paul Rayson. Boss: Bayesian optimization over string spaces. *Advances in neural information processing systems*, 33:15476–15486, 2020.
- AkshatKumar Nigam, Pascal Friederich, Mario Krenn, and Alan Aspuru-Guzik. Augmenting genetic algorithms with deep neural networks for exploring the chemical space. In *International Conference on Learning Representations*, 2020.
- Pascal Notin, José Miguel Hernández-Lobato, and Yarin Gal. Improving black-box optimization in vae latent space using decoder uncertainty. *Advances in Neural Information Processing Systems*, 34, 2021.
- Gregor Simm, Robert Pinsler, and José Miguel Hernández-Lobato. Reinforcement learning for molecular design guided by quantum mechanics. In *International Conference on Machine Learning*, pp. 8959–8969. PMLR, 2020.
- Gregor N. C. Simm, Robert Pinsler, Gábor Csányi, and José Miguel Hernández-Lobato. Symmetry-aware actor-critic for 3d molecular design. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=jEYKjPE1xYN>.
- Austin Tripp, Erik Daxberger, and José Miguel Hernández-Lobato. Sample-efficient optimization in the latent space of deep generative models via weighted retraining. *Advances in Neural Information Processing Systems*, 33:11259–11272, 2020.
- Austin Tripp, Gregor NC Simm, and José Miguel Hernández-Lobato. A fresh look at de novo molecular design benchmarks. In *NeurIPS 2021 AI for Science Workshop*, 2021.
- Oleg Trott and Arthur J Olson. Autodock vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *Journal of computational chemistry*, 31(2):455–461, 2010.
- Scott A Wildman and Gordon M Crippen. Prediction of physicochemical parameters by atomic contributions. *J. Chem. Inf. Comput. Sci.*, 39(5):868–873, 1999.
- Yutong Xie, Chence Shi, Hao Zhou, Yuwei Yang, Weinan Zhang, Yong Yu, and Lei Li. {MARS}: Markov molecular sampling for multi-objective drug discovery. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=kHSu4ebxFXy>.
- Chencheng Xu, Qiao Liu, Minlie Huang, and Tao Jiang. Reinforced molecular optimization with neighborhood-controlled grammars. *Advances in Neural Information Processing Systems*, 33: 8366–8377, 2020.
- Kevin Yang, Wengong Jin, Kyle Swanson, Regina Barzilay, and Tommi Jaakkola. Improving molecular design by stochastic iterative target augmentation. In *International Conference on Machine Learning*, pp. 10716–10726. PMLR, 2020.
- Soojung Yang, Doyeong Hwang, Seul Lee, Seongok Ryu, and Sung Ju Hwang. Hit and lead discovery with explorative rl and fragment-based molecule generation. *Advances in Neural Information Processing Systems*, 34, 2021.
- Jiaxuan You, Bowen Liu, Rex Ying, Vijay Pande, and Jure Leskovec. Graph convolutional policy network for goal-directed molecular graph generation. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pp. 6412–6422, 2018.



## A RECENT WORKS IN MACHINE LEARNING FOCUSING ON MOLECULAR OPTIMIZATION

The following is a non-exhaustive list of papers doing molecular optimization from three large machine learning conferences.

**ICLR** (Nigam et al., 2020; Xie et al., 2021; Simm et al., 2021; Ahn et al., 2022; Fu et al., 2022; Maziarz et al., 2022; Gao et al., 2022)

**ICML** (Kusner et al., 2017; Jin et al., 2018; Liu et al., 2020; Yang et al., 2020; Jin et al., 2020b;a; Simm et al., 2020; Luo et al., 2021)

**NeurIPS** (Liu et al., 2018; You et al., 2018; Bradshaw et al., 2019; Chenthamarakshan et al., 2020; Ahn et al., 2020; Moss et al., 2020; Tripp et al., 2020; Bradshaw et al., 2020; Xu et al., 2020; Mollaysa et al., 2020; Bengio et al., 2021; Yang et al., 2021; Notin et al., 2021)