

THESIS FOR THE DEGREE OF LICENTIATE OF ENGINEERING

Sampling from molecular unnormalized distributions with Deep Generative Models

*Toward the acceleration of molecular design and conformational
sampling with Deep Learning*

JUAN VIGUERA DIEZ

Department of Computer Science and Engineering
CHALMERS UNIVERSITY OF TECHNOLOGY | UNIVERSITY OF GOTHENBURG
Gothenburg, Sweden, 2024

Sampling from molecular unnormalized distributions with Deep Generative Models

Toward the acceleration of molecular design and conformational sampling with Deep Learning

JUAN VIGUERA DIEZ

© Juan Viguera Diez, 2024
except where otherwise stated.
All rights reserved.

ISSN 1652-876X

Department of Computer Science and Engineering
Division of Data Science and AI
Simon Olsson's group
Chalmers University of Technology | University of Gothenburg
SE-412 96 Göteborg,
Sweden
Phone: +46(0)31 772 1000

Printed by Chalmers Digitaltryck,
Gothenburg, Sweden 2024.

To my family
A mi familia

Sampling from molecular unnormalized distributions with Deep Generative Models

Toward the acceleration of molecular design and conformational sampling with Deep Learning

JUAN VIGUERA DIEZ

*Department of Computer Science and Engineering
Chalmers University of Technology | University of Gothenburg*

Abstract

This thesis investigates how Deep Generative Models (DGMs) can address important drug discovery problems involving sampling from unnormalized distributions. It consists of two papers focusing on this challenge's aspects: molecular design and conformational sampling. The first paper proposes a new training scheme to fine-tune graph-based DGMs for *de novo* molecular design. Our method can produce molecules with specific properties even when they are scarce or missing in the training data and outperforms previously reported graph-based methods on predicted dopamine receptor type D2 activity while maintaining diversity. The second paper develops Surrogate Model-Assisted Molecular Dynamics (SMA-MD), which combines a DGM with statistical re-weighting and short Molecular Dynamics simulations to generate equilibrium ensembles of molecules. SMA-MD can produce more diverse and lower energy ensembles than conventional molecular dynamics simulations. These contributions constitute important stepping stones towards the automation of the drug discovery process.

Keywords

Cheminformatics, machine learning, drug discovery, generative models, conformational sampling, Boltzmann generators

List of Publications

Appended publications

This thesis is based on the following publications:

- [**Paper I**] S. Romeo Atance, **J. Viguera Diez**, O. Engkvist, S. Olsson, R. Mercado, *De Novo Drug Design Using Reinforcement Learning with Graph-Based Deep Generative Models*
J. Chem. Inf. Model. 2022, 62, 20, 4863–4872.
doi.org/10.1021/acs.jcim.2c00838
- [**Paper II**] **J. Viguera Diez**, S. Romeo Atance, O. Engkvist, S. Olsson, *Generation of conformational ensembles of small molecules via Surrogate Model-Assisted Molecular Dynamics*
Submitted, under revision.
doi.org/10.26434/chemrxiv-2023-sx61w

Other publications

The following publications were published during my PhD studies, or are currently in submission/under revision. However, they are not appended to this thesis, due to contents overlapping that of appended publications or contents not related to the thesis.

- [a] **J. Viguera Diez**, S. Romeo Atance, O. Engkvist, R. Mercado, S. Olsson, *A transferable Boltzmann generator for small-molecules conformers* *ELLIS Machine Learning for Molecule Discovery Workshop (November 2021)*, 7.

Acknowledgment

First, I would like to thank my academic supervisor and mentor, Simon Olsson. Thank you for your guidance, support, and the freedom you have given me during these years. You are an inspiration and an example to follow for me. I would also like to thank my industrial advisor, Ola Engkvist. Thank you for your support, flexibility, and expertise. It is a pleasure to learn from your scientific knowledge and experience. I am also grateful to the other members of the committee, Devdatt Dubhashi and Morteza Haghiri Chehreghani, for their support and valuable feedback on my research. I would like to thank Søren Hauberg for accepting to be the discussion leader for my Licentiate defense.

I thank my colleagues and friends, both at the Data Science and AI (DSAI) division at Chalmers and the Molecular AI team at AstraZeneca, for their camaraderie and for sharing their experiences with me. I have learned so much from you and I am really happy to share this journey with you. I especially thank Rocío Mercado for her invaluable support and expertise during her time at AstraZeneca and now at Chalmers. At DSAI, I would like to especially thank Chris, Mathias, Janosch, Tobias, Télió, Lena, Lovisa, Riccardo, Anton, Adam, Newton, and Jon. Thanks to Claes Andersson for helping me to develop my teaching skills. At Molecular AI, I would like to especially thank Alessandro, Jon Paul, Samuel, Marco, Yasmine, Gökçe, Thomas, Tomas, Jiazhen, Annie, Helen, Peter, Emma, Rosa, Varvara, Paula and Vincenzo. I thank Atanas Patronov for his supervision over the first months of my PhD. I have the pleasure of being colleagues, both at Chalmers and AstraZeneca, with Simon, Hampus, Ross, and Emma and I am thankful for sharing this experience with them.

Last, but certainly not least, I thank my friends and family. I especially thank my partner, Sara, for her support, and contributions. I am really thankful to have you in my life. I am forever thankful to my parents Marta and Luis and my brother Rubén for their unconditional love and support. Thank you for always being there and believing in me. Nothing of this would have been possible without all of you. ¡Gracias!

This work was partially supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP), funded by the Knut and Alice Wallenberg Foundation. I am grateful for the funding and the educational and networking opportunities provided by WASP.

Contents

Abstract	iii
List of Publications	v
Acknowledgement	vii
I Introductory chapters	1
1 Introduction	3
2 Background	5
2.1 The drug discovery process	5
2.2 Molecular representations	6
2.2.1 Symmetries in molecular Machine Learning	8
2.3 Machine Learning accelerated drug design	9
2.3.1 <i>De novo</i> design	10
2.3.2 Molecular property prediction	10
2.3.3 Conformational sampling	11
3 Summary of Included Papers	13
3.1 <i>De Novo</i> Drug Design Using Reinforcement Learning with Graph-Based Deep Generative Models	13
3.2 Generation of conformational ensembles of small molecules via Surrogate Model-Assisted Molecular Dynamics	15
4 Discussion and Future Work	17
Bibliography	19
II Appended Papers	27
Paper I - <i>De Novo</i> Drug Design Using Reinforcement Learning with Graph-Based Deep Generative Models	

Paper II - Generation of conformational ensembles of small molecules via Surrogate Model-Assisted Molecular Dynamics

Part I

Introductory chapters

Chapter 1

Introduction

Drug discovery is the process of identifying and developing new compounds that can modulate the activity of a biological target, such as a protein or a gene, for therapeutic purposes. Drug discovery is a complex, costly, and time-consuming endeavor, requiring multidisciplinary expertise and collaboration. According to some estimates, it takes an average of 10 years and 2.6 billion dollars to bring a new drug to the market [1]. Hence, finding ways to speed up and simplify this process is crucial.

Machine Learning (ML) technologies can streamline and expedite various aspects of drug design. They can help identify and validate new targets by analyzing massive amounts of biological and genomic data [2], [3], enhance the design and optimization of drug candidates by using generative models and molecular simulations [4]–[6], and improve the prediction of drug properties and safety [7]–[9]. In addition, natural language processing and computer vision techniques can facilitate the clinical development and testing of drugs [3], [10].

Sampling from unnormalized distributions is an important open scientific problem with applications in drug discovery [11], material science [12] and machine learning [13], [14]. The objective is to generate samples, x , following a certain probability density $p(x)$, proportional to some function from which the normalizing constant is unknown. This function is usually an exponential of some fitness function $f(x)$. There are several approaches to sample from unnormalized distributions, such as simulation-based methods [15], [16] or Markov Chain Monte Carlo [17]–[19]. However, they perform poorly for high-dimensional and topologically complex data, such as molecules. Deep Generative Models (DGMs) are a class of machine learning models that can be used to generate new data samples that are similar to the training data. DGMs are topologically flexible and well-behaved in high-dimensional spaces.

This thesis explores two solutions based on DGMs to sample from unnormalized molecular distributions $p(x)$. First, we explore molecular design, in which x is a molecular graph and f is an oracle function assessing the desirability of a molecule. Second, we delve into conformational sampling, where x will be a 3-dimensional arrangement of atoms and f will be the negative reduced energy.

The first paper [20] presents a novel training scheme to fine-tune graph-

based Deep Generative Models for *de novo* molecular design. This scheme guides the model to generate molecules with desired properties, even when they are rare or absent in the training data. We use a Graph Neural Network to model the action probability distributions for building molecular graphs, and introduce a memory-aware loss function to speed up and stabilize learning. We demonstrate the effectiveness of this approach on several design tasks, especially for generating molecules with predicted dopamine receptor type D2 activity.

The second paper [21] introduces a new method for generating equilibrium ensembles of molecules that combines a DGM with statistical reweighting and short Molecular Dynamics (MD) simulations. The method, called Surrogate Model-Assisted Molecular Dynamics (SMA-MD), can produce more diverse and lower energy ensembles than conventional MD simulations, and can also estimate implicit solvation free energies. SMA-MD is demonstrated to be an efficient and transferable approach for sampling from the Boltzmann distribution of small molecular systems.

Finally, we will discuss how these methods could impact drug design pipelines in the future and what important problems towards the automation of drug design remain unsolved.

Chapter 2

Background

2.1 The drug discovery process

Drug discovery is the process of finding new medications based on the knowledge of a biological target. A biological target is a key molecule involved in a particular metabolic or signaling pathway that is associated with a specific disease condition or pathology. Drug discovery is a multi-step process that aims to create molecules that can bind to the target and modulate its function, resulting in a therapeutic benefit to the patient [22] (Figure 2.1).

The first stage of drug discovery involves target identification and validation. Target identification [23] is the process of finding a biological target, such as a protein or a gene, that is involved in the disease mechanism and can be modulated by a drug. Popular techniques include genomics [24] and proteomics [25], which aim to describe the structure, function, and dynamics of genes/proteins in the human body. Target validation is the process of confirming that the target is indeed relevant to the disease and is performed using methods such as gene knockouts. Gene knockouts [26] consist of removing or inactivating specific genes within an organism's genome to determine the effect on the disease mechanism.

The second stage is drug design and consists of lead discovery and optimization. Lead discovery is the process of screening large libraries of chemical compounds for those that have potential activity against the target. One important method is *de novo* design [27], which is a computational technique to generate novel molecular structures that have desired properties or functions. Lead optimization is the process of improving the properties of lead compounds, such as potency, selectivity, solubility, stability, and toxicity, to make them suitable for further development. One key technique is Quantitative Structure-Activity Relationship (QSAR) modeling [28], [29], which involves predicting the biological activity of a compound based on its chemical structure.

Finally, drug development covers the steps taken to convert the lead compound into an approved drug product for human use. Drug development includes preclinical research to evaluate the safety and efficacy of the compound in animal models, clinical research to test its effects in human volunteers,

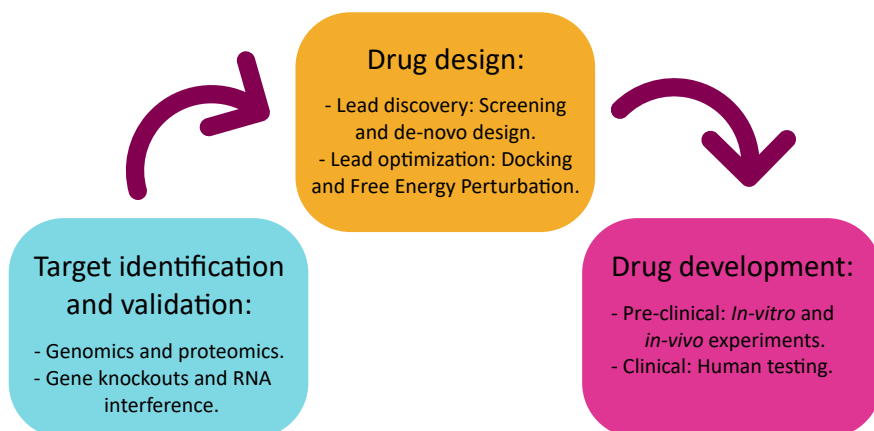


Figure 2.1: The drug discovery process: stages and key techniques

regulatory review to obtain approval from authorities such as the European Medicines Agency or the Food and Drug Administration, and post-market monitoring to ensure its safety and quality after launch.

2.2 Molecular representations

In this work, we restrict ourselves to representing molecules at the Born-Oppenheimer level of detail, which specifies the position of the nuclei of the atoms in a molecule. Molecules are dynamic and adopt a variety of 3-dimensional structures called conformations. In equilibrium, the ensemble of conformations, \mathbf{x} , that a molecule can adopt follows the Boltzmann distribution,

$$\mu(\mathbf{x}) = \mathcal{Z}^{-1} \exp(-\beta U(\mathbf{x})), \text{ with } \mathcal{Z} = \int d\mathbf{x} \exp(-\beta U(\mathbf{x})), \quad (2.1)$$

where $U(\mathbf{x})$ is the potential energy and β is the inverse temperature. Conformational ensembles contain substantial information about molecules, such as their structural diversity or relative populations. However, generating representative conformational ensembles is challenging due to the long time and length scales involved in transitions among conformations.

In some situations, it is possible to approximate the Boltzmann distribution by only accounting for the local maxima of the distribution. This set of representative conformations is referred to as conformers. The conformer with minimum energy is referred to as the ground state and is often the most statistically representative. Conformers are simpler yet less complete descriptions of molecules when compared with conformational ensembles.

Conformational ensembles and conformers can be encoded as arrays of 3-dimensional structures. One simple way to store this information is to use the Cartesian Coordinates (CCs) of the atoms as shown in Figure 2.2 (a). Nevertheless, this representation is sensitive to global roto-translations of the

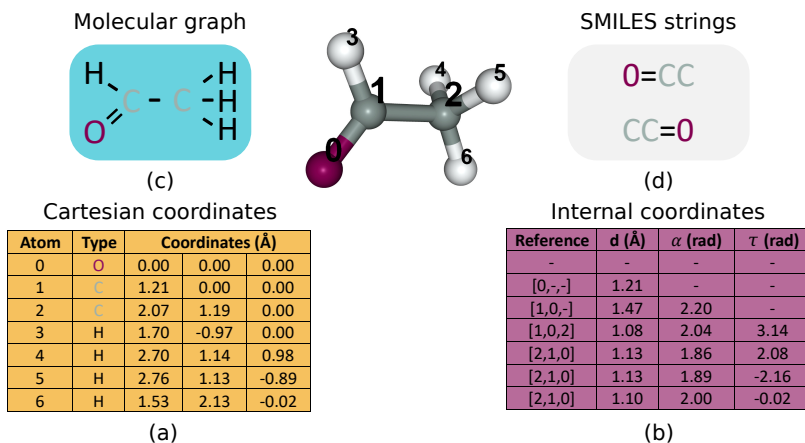


Figure 2.2: Molecular representations for acetaldehyde: Cartesian Coordinates (a), Internal Coordinates (b), consisting of distances (d), angles (α), and dihedral angles (τ), molecular graph (c), and SMILES strings (d).

structure, which is not desirable in many applications. For this reason, Internal Coordinates (ICs) may be used instead. ICs are a combination of bond lengths, bond angles, and dihedral angles which can be converted to CCs, and are illustrated in Figure 2.2 (b). The collection of ICs describing a complete structure, along with atom types and reference atoms is known as the Z-matrix.

For some applications, a simpler molecular representation including solely topological information may be preferred. The connectivity of a molecule can be represented using its chemical graph, which is a collection of nodes (atoms) and edges (chemical bonds). An example is depicted in Figure 2.2 (c). The information contained in a molecular graph can be stored in a string of characters using the Simplified Molecular-Input Line-Entry System (SMILES) [30] as exemplified in Figure 2.2 (d). SMILES is widely used for exchanging and storing molecular structures in databases and software. A molecular graph can be inferred from a SMILES string, but several SMILES strings can correspond to a given chemical graph. This can cause issues for some models, as they may generate inconsistent outputs for the same molecule depending on the input SMILES.

The optimal molecular representation depends on the problem at hand and the available computational resources. For example, if the goal is to accurately predict the biological activity of a molecule, then a representation that captures the 3-dimensional structure of the molecule may be required. On the other hand, if the goal is to perform a large-scale virtual screening of a database of compounds, then the molecular graph may be sufficient. More elaborated molecular descriptors provide more accurate predictions at the cost of increasing compute time. Properties of different molecular representations are summarized in Table 2.1.

	Advantages	Drawbacks
CCs	Simple	Sensitive to roto-translations
ICs	Invariant to roto-translations	Requires specification of reference atoms
Graph	Unique	Lack of 3-dimensional information
SMILES	Simple and lightweight	Not unique. Lack of 3-dimensional information.

Table 2.1: Advantages and drawbacks of different molecular representations: Cartesian Coordinates (CCs), Internal Coordinates (ICs), chemical graphs, and Simplified Molecular-Input Line-Entry System (SMILES) string.

2.2.1 Symmetries in molecular Machine Learning

Symmetries are an important consideration when designing Machine Learning solutions for chemistry [31]–[34]. For example, if we were to predict the forces on the atoms of an isolated molecule, it is desirable for the prediction to rotate along with the input molecule. However, if we were to predict its strain energy, the prediction should not change with global roto-translations of the input. We say that the forces model should be equivariant w.r.t. rotations while the energy model should be invariant. Similarly, when we use graph-based representations, we want our model to be invariant to the ordering of the atoms, and different SMILES strings for the same chemical graph. In general, given a model m and a transformation t , we say that m is equivariant w.r.t. t if

$$m(t(x)) = t(m(x)) \quad (2.2)$$

and invariant if

$$m(t(x)) = m(x). \quad (2.3)$$

Considering symmetries is important in practice since it reduces sample complexity (makes training more efficient and alleviates the need for data augmentation) and reduces the hypothesis space of learnable models [32], [35].

Graph Neural Networks (GNNs) [36], [37] are a widely used method for incorporating symmetries in molecular ML models. GNNs consist of layers that update the node features by aggregating information from their neighbors. This operation is known as message passing, and it can be expressed as:

$$x_i^{(l+1)} = \phi \left(x_i^{(l)}, \bigoplus_{j \in \mathcal{N}(i)} \psi \left(x_i^{(l)}, x_j^{(l)} \right) \right), \quad (2.4)$$

where the feature vector $x_i^{(l)}$ of atom (node) i in layer l gets updated combining the messages $\psi \left(x_i^{(l)}, x_j^{(l)} \right)$ from the neighboring atoms $\mathcal{N}(i)$ using an aggregation function \bigoplus , and an update function ϕ . By using this message passing scheme, GNNs are invariant to the permutation of nodes in the graph (atom order). Moreover, GNNs can also incorporate other symmetries, such as rotation or translation, by designing the functions ϕ and ψ to be either invariant or equivariant to the desired transformations. For example, invariance

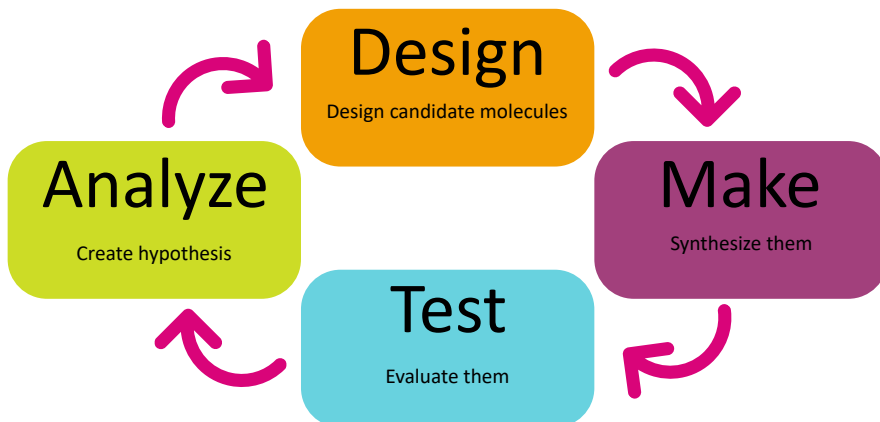


Figure 2.3: Design-Make-Test-Analyze cycle. A model for drug design.

and equivariance to rotations can be achieved, respectively, by incorporating interatomic distances and relative positions as geometric features.

2.3 Machine Learning accelerated drug design

Drug design is a stage of drug discovery involving iterative trial-and-error testing. For example, how does property A change when group B is changed in this molecule? This situation is often represented with a model known as the Design-Make-Test-Analyze (DMTA) cycle, illustrated in Figure 2.3.

The DMTA cycle is often time-consuming and costly, but ML can offer a promising solution for its automation and acceleration. In the design phase, DGMs can be used to generate novel molecules with desired features such as high binding affinity or low toxicity [4], [38]–[40]. In the make phase, better synthesis routes may be found more efficiently by using automated methods [9], [41]–[43]. In the test phase, ML can enable faster predictions of important attributes of potential drug candidates, such as solubility, stability, or pharmacokinetics [44], [45]. For example, one can use regression or classification models to estimate the properties of new molecules based on their structural or physicochemical features. Another important example is using DGMs for generating 3-dimensional conformations of molecules from which properties may be estimated [5], [31], [46]–[52]. In the analyze phase, ML can generate plausible interpretations of the test phase results and provide insights into the structure-activity relationship or mechanism of action of new molecules [8], [53], [54]. Moreover, new directions for further optimization or discovery may be suggested by using techniques such as Active Learning or Reinforcement Learning [55], [56].

2.3.1 *De novo* design

De novo design is a computational approach to generate novel molecular structures that have desired properties or functions. This tool is used during lead generation and optimization stages. DGMs are increasingly being adopted for *de novo* design, leading to the emerging field of generative chemistry. Most approaches start from training a prior generative model on large chemical datasets, which can be seen as ‘foundation models’ for chemistry. However, these models are usually expensive to train and may not generate molecules that satisfy the design requirements. Therefore, transfer learning [57], [58] is used to bias the prior model towards regions of chemical space that are more desirable. There are two main approaches. The first one consists of fine-tuning the prior model on a smaller dataset of molecules satisfying the design constraints [59]. The second one relies on the availability of a scoring model, assessing the desirability of molecules, and consists of biasing the prior model to promote the generation of highly scored compounds [4]. If target data are scarce, usually the first approach is preferred, but if a reliable model is available, it may be exploited by following the second approach.

One popular mathematical formulation of the second scenario consists of sampling from an unnormalized distribution with support over the compositional space of molecules. The objective is to sample from the distribution,

$$p(x) \propto \exp(\sigma S(x) + \log p_{\mathbb{P}}(x)), \quad (2.5)$$

where x is a molecular graph, S is a scoring model, $p_{\mathbb{P}}(x)$ is the probability density of finding molecule x sampling from an pre-trained unbiased model and σ is a parameter balancing both contributions. Intuitively, the first term modulates the probability of the pre-trained model so that high-scoring compounds are sampled more likely compared to those scored poorly. Computing the normalizing constant of this distribution is not tractable as it requires accounting for every possible molecular graph.

Several types of generative models have been proposed to model the distribution $p(x)$ [60]. Previous work based on SMILES strings has used Recurrent Neural Networks (RNNs) [4], [59], [61], Variational Autoencoders (VAEs) [62], [63], or Generative Adversarial Networks (GANs) [64], [65]. Methods using a graph representation have used different types of GNNs [66] such as Gated Graph Neural Networks (GGNN) [38] or Graph Convolutional Neural Networks (GCNNs) [67] to iteratively sample actions that build up a chemical graph. More recently, 3D-generative models [39], [68] have explored molecular generation directly in the binding pocket, potentially accounting explicitly for physical interactions. Prevalent methods are diffusion models [69]–[71] and normalizing flows [72] powered by Equivariant Neural Networks [31], [32].

2.3.2 Molecular property prediction

Predicting the properties of molecules is an essential part of drug design since it allows for the identification of promising candidates and reduces the cost and time of experimental testing. Naturally, molecular properties are a core

component of scoring functions. Here, we focus on two methods for predicting the properties of molecules, graph-based and conformational ensemble-based methods.

On the one hand, graph-based models for property prediction take the molecular graph (or SMILES) as input and generate predictions as outputs. Diverse models can be used in this context such as logistic regression, support vector machines, random forests or neural networks [73]. One particularly important type of models are QSAR models [7], which are used to predict the activity of drug candidates w.r.t. a given target. QSAR models can be classified into regression models, which predict continuous activity values, or classification models, which predict categorical activity values [8]. These models tend to be lightweight and fast and therefore are often used for building scoring functions guiding generative models.

On the other hand, conformational ensemble-based methods rely on inferring molecular properties from a representative conformational ensemble, $\{\mathbf{x}_i\}$. Given independent and identically distributed (i.i.d.) conformations sampled from the Boltzmann distribution, a molecular property, O can be computed via the Monte Carlo estimator:

$$O = \mathbb{E}_{\mathbf{x} \sim \mu(\mathbf{x})}[o(\mathbf{x})] \approx \frac{1}{N} \sum_{i=1}^N o(\mathbf{x}_i), \quad \mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \mu(\mathbf{x}), \quad (2.6)$$

where $o(\mathbf{x})$ computes the microscopic contribution of a property for a conformation. For example, if O is the distance between two atoms, $o(\mathbf{x})$ would be the distance between those atoms in each conformation. Although conformational ensemble-based methods tend to be more accurate than graph-based, this comes at the cost of generating representative ensembles.

2.3.3 Conformational sampling

Conformational ensemble-based methods rely on i.i.d. samples from the Boltzmann distribution, Equation 2.1. Therefore, generating these samples, referred to as conformational sampling, is a crucial step. In general, computing the normalizing constant of the Boltzmann distribution, the partition function \mathcal{Z} , is not tractable. Therefore, conformational sampling consists of sampling from an unnormalized distribution with support over the conformational space.

Key solutions to conformational sampling such as molecular dynamics simulations to Markov-Chain Monte Carlo (MCMC) methods require many simulation steps to generate representative ensembles, especially for high-dimensional and meta-stable systems. A promising approach to this problem consists of approximating the Boltzmann distribution of a molecular system with a DGM. To be an effective solution, the DGM must allow efficient sampling and exact likelihood evaluation. Efficient i.i.d. sampling allows to side-step iterative simulation methods, and exact likelihood evaluation allows to recover unbiased samples through importance sampling. Methods implementing this solution are called Boltzmann Generators (BGs) [5], [46], [48]. However, currently, BGs suffer from limited transferability across molecular systems and conformational sampling remains an important open problem.

Chapter 3

Summary of Included Papers

3.1 *De Novo* Drug Design Using Reinforcement Learning with Graph-Based Deep Generative Models

In the first article, we propose a new training scheme to fine-tune graph-based DGMs for *de novo* molecular design tasks. We show how our computational framework can successfully guide a pre-trained generative model toward the generation of molecules with a specific property profile, even when such molecules are not present in the training set and unlikely to be generated by the pre-trained model. We explored the following tasks: generating molecules of decreasing/increasing size, increasing drug-likeness, and increasing bioactivity.

We use GraphINVENT [38] as a graph-based molecular DGM. GraphINVENT is based on a Gated Graph Neural Network (GGNN) that generates molecules by iteratively sampling actions that build upon an input graph. The action space is divided into three possible actions: add atom, add bond, and terminate graph. The model is trained by minimizing the Kullback-Leibler divergence between target and predicted action probability distributions (APDs).

Our Reinforcement Learning framework uses a memory-aware loss that keeps track of the best agent so far and is updated every few learning steps. By doing so, we remind the current agent of sets of actions that can lead to high-scoring compounds, in turn accelerating and improving agent learning. The scoring model is designed for each specific optimization task and can be based on simple rules or more complex models such as QSAR models.

We tested our framework by fine-tuning a pre-trained graph-based DGM to favor property profiles relevant to drug design, including increasing pharmacological activity. We model bioactivity using a QSAR model for dopamine receptor type D2 (DRD2) activity. Optimization for DRD2 activity is a widely used *de novo* design bioactivity benchmark and allows us to easily compare to

previous work. We achieve models that generate diverse compounds with predicted DRD2 activity for 97 % of sampled molecules, outperforming previously reported graph-based methods on this metric.

Our contribution is an important stepping stone toward the design of more advanced molecular DGMs which will allow scientists to efficiently traverse the chemical space in search of promising molecules. We believe the use of DGMs in fields such as drug design has the potential to help chemists come up with new ideas and accelerate the complex process of molecular discovery.

3.2 Generation of conformational ensembles of small molecules via Surrogate Model-Assisted Molecular Dynamics

In the second paper, we present a new method for generating equilibrium conformational ensembles of molecules, called Surrogate Model-Assisted Molecular Dynamics (SMA-MD). This method combines a DGM that samples the slow degrees of freedom of molecules with a reweighting and short simulation step that equilibrates the fast degrees of freedom. SMA-MD can generate more diverse and physically realistic ensembles than conventional MD simulations.

We use a two-step procedure to generate molecular conformations. First, we use a deterministic algorithm to generate the local structure of each atom, and then we use a diffusion model to sample the torsion angles of rotatable bonds [52]. The diffusion model is trained on MD simulations of small non-cyclic molecules. Second, we reweight the generated conformations against the Boltzmann distribution and run short parallel MD simulations to thermalize and mix the fast degrees of freedom.

We evaluate our method by comparing it with MD and Replica Exchange (RE) simulations on various metrics, such as conformer generation, potential energy, free energy of solvation, and slow transitions. We show that SMA-MD outperforms MD in generating more diverse and energetically favorable ensembles, and matches RE in capturing the relevant states and properties of molecules.

We conclude that SMA-MD is an efficient and robust method for sampling from the Boltzmann distribution of molecules. We highlight the advantages of SMA-MD over MD, such as data aggregation, parallelization, and independence of initial conditions. We also discuss the limitations and future directions of SMA-MD, such as extending it to cyclic molecules, improving the computational cost of sampling, and training Boltzmann surrogates with large-scale data.

SMA-MD shows promising results toward accelerating the generation of representative conformational ensembles of molecules with DGMs. As such, SMA-MD is a step toward faster methods for predicting molecular properties, which is fundamental in drug design.

Chapter 4

Discussion and Future Work

This thesis investigates how DGMs can address important drug discovery problems involving sampling from unnormalized distributions. It includes two papers focusing on different aspects of this challenge: molecular design and conformational sampling. In the first paper, we proposed a training scheme to fine-tune graph-based DGMs for *de novo* molecular design, which can generate molecules with specific properties even when they are scarce or missing in the training data. In the second paper, we developed Surrogate Model-Assisted Molecular Dynamics (SMA-MD), which combines a DGM with statistical reweighting and short MD simulations to generate equilibrium ensembles of molecules. SMA-MD can produce more diverse and lower energy ensembles than conventional MD simulations.

These contributions constitute important stepping stones towards the automation of the drug discovery process. They demonstrate the potential of DGMs to sample from complex and high-dimensional molecular spaces, and to optimize molecules for multiple criteria. They also highlight the challenges and limitations of these models.

We demonstrated the effectiveness of our approach for fine-tuning graph-based DGMs by generating 97 % active molecules for the dopamine receptor type D2. However, the QSAR model that we used for both fine-tuning the DGM and evaluating the method only relied on graph-level information. A more challenging and realistic task is to predict pharmacological activity while taking into account 3-dimensional information. Moreover, *de novo* design is inherently a 3-dimensional problem, and therefore methods that can generate molecules directly in 3D and optimize them for the target interaction are a promising direction for future work.

DGMs are a powerful tool for conformational sampling and we have demonstrated that methods based on DGMS can outperform classic MD simulations in the context of small molecules. However, transferability across molecular systems is challenging and it is unclear if our conclusions generalize to more complex systems such as drug-like molecules or proteins. This is an interesting

avenue for future work. Additionally, the models presented in this thesis cannot provide information about the kinetics of molecules, which is necessary for predicting some properties. Building transferable DGMs of the transition probability of molecular conformations remains an open problem and will also be explored in future research. Finally, exploring how conformational DGMs can speed up binding affinity predictions remains an important open problem that we may explore in the future.

Bibliography

- [1] J. A. DiMasi, H. G. Grabowski and R. W. Hansen, “Innovation in the pharmaceutical industry: New estimates of r&d costs,” *Journal of Health Economics*, vol. 47, pp. 20–33, 2016, ISSN: 0167-6296. DOI: <https://doi.org/10.1016/j.jhealeco.2016.01.012>. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0167629616000291> (cit. on p. 3).
- [2] F. W. Pun, I. V. Ozerov and A. Zhavoronkov, “Ai-powered therapeutic target discovery,” *Trends in Pharmacological Sciences*, vol. 44, no. 9, pp. 561–572, 2023, ISSN: 0165-6147. DOI: <https://doi.org/10.1016/j.tips.2023.06.010>. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0165614723001372> (cit. on p. 3).
- [3] Y. You, X. Lai, Y. Pan *et al.*, “Artificial intelligence in cancer target identification and drug discovery,” *Signal Transduction and Targeted Therapy*, vol. 7, no. 1, May 2022, ISSN: 2059-3635. DOI: 10.1038/s41392-022-00994-0. [Online]. Available: <http://dx.doi.org/10.1038/s41392-022-00994-0> (cit. on p. 3).
- [4] M. Olivecrona, T. Blaschke, O. Engkvist and H. Chen, “Molecular de-novo design through deep reinforcement learning,” *Journal of Cheminformatics*, vol. 9, no. 1, Sep. 2017, ISSN: 1758-2946. DOI: 10.1186/s13321-017-0235-x. [Online]. Available: <http://dx.doi.org/10.1186/s13321-017-0235-x> (cit. on pp. 3, 9, 10).
- [5] F. Noé, S. Olsson, J. Köhler and H. Wu, *Boltzmann generators – sampling equilibrium states of many-body systems with deep learning*, 2019. arXiv: 1812.01729 [stat.ML] (cit. on pp. 3, 9, 11).
- [6] S. Chmiela, V. Vassilev-Galindo, O. T. Unke *et al.*, “Accurate global machine learning force fields for molecules with hundreds of atoms,” *Science Advances*, vol. 9, no. 2, Jan. 2023, ISSN: 2375-2548. DOI: 10.1126/sciadv.adf0873. [Online]. Available: <http://dx.doi.org/10.1126/sciadv.adf0873> (cit. on p. 3).
- [7] P.-C. Kotsias, J. Arús-Pous, H. Chen, O. Engkvist, C. Tyrchan and E. J. Bjerrum, “Direct steering of de novo molecular generation using descriptor conditional recurrent neural networks (crnns),” Nov. 2019. DOI: 10.26434/chemrxiv.9860906.v2. [Online]. Available: <http://dx.doi.org/10.26434/chemrxiv.9860906.v2> (cit. on pp. 3, 11).

- [8] M. R. Keyvanpour and M. B. Shirzad, “An analysis of QSAR research based on machine learning concepts,” en, *Curr Drug Discov Technol*, vol. 18, no. 1, pp. 17–30, 2021 (cit. on pp. 3, 9, 11).
- [9] S. Genheden, A. Thakkar, V. Chadimová, J.-L. Reymond, O. Engkvist and E. Bjerrum, “Aizynthfinder: A fast, robust and flexible open-source software for retrosynthetic planning,” *Journal of Cheminformatics*, vol. 12, no. 1, Nov. 2020, ISSN: 1758-2946. DOI: 10.1186/s13321-020-00472-1. [Online]. Available: <http://dx.doi.org/10.1186/s13321-020-00472-1> (cit. on pp. 3, 9).
- [10] R. Qureshi, M. Irfan, T. M. Gondal *et al.*, “Ai in drug discovery and its clinical relevance,” *Heliyon*, vol. 9, no. 7, e17575, Jul. 2023, ISSN: 2405-8440. DOI: 10.1016/j.heliyon.2023.e17575. [Online]. Available: <http://dx.doi.org/10.1016/j.heliyon.2023.e17575> (cit. on p. 3).
- [11] Z. Cournia, C. Chipot, B. Roux, D. M. York and W. Sherman, “Free energy methods in drug discovery—introduction,” in *ACS Symposium Series*. American Chemical Society, Nov. 2021, 1–38, ISBN: 9780841298057. DOI: 10.1021/bk-2021-1397.ch001. [Online]. Available: <http://dx.doi.org/10.1021/bk-2021-1397.ch001> (cit. on p. 3).
- [12] J. Köhler, M. Invernizzi, P. de Haan and F. Noé, *Rigid body flows for sampling molecular crystal structures*, 2023. arXiv: 2301.11355 [cs.LG] (cit. on p. 3).
- [13] Y. W. Teh, M. Welling, S. Osindero and G. E. Hinton, “Energy-based models for sparse overcomplete representations,” *J. Mach. Learn. Res.*, vol. 4, no. null, 1235–1260, Dec. 2003, ISSN: 1532-4435 (cit. on p. 3).
- [14] T. Haarnoja, H. Tang, P. Abbeel and S. Levine, *Reinforcement learning with deep energy-based policies*, 2017. arXiv: 1702.08165 [cs.LG] (cit. on p. 3).
- [15] S. A. Hollingsworth and R. O. Dror, “Molecular dynamics simulation for all,” *Neuron*, vol. 99, no. 6, 1129–1143, Sep. 2018, ISSN: 0896-6273. DOI: 10.1016/j.neuron.2018.08.011. [Online]. Available: <http://dx.doi.org/10.1016/j.neuron.2018.08.011> (cit. on p. 3).
- [16] J. D. Durrant and J. A. McCammon, “Molecular dynamics simulations and drug discovery,” *BMC Biology*, vol. 9, no. 1, Oct. 2011, ISSN: 1741-7007. DOI: 10.1186/1741-7007-9-71. [Online]. Available: <http://dx.doi.org/10.1186/1741-7007-9-71> (cit. on p. 3).
- [17] D. van Ravenzwaaij, P. Cassey and S. D. Brown, “A simple introduction to markov chain monte-carlo sampling,” *Psychonomic Bulletin & Review*, vol. 25, no. 1, 143–154, Mar. 2016, ISSN: 1531-5320. DOI: 10.3758/s13423-016-1015-8. [Online]. Available: <http://dx.doi.org/10.3758/s13423-016-1015-8> (cit. on p. 3).

- [18] G. Casella, C. P. Robert and M. T. Wells, “Generalized accept-reject sampling schemes,” in *A Festschrift for Herman Rubin*. Institute of Mathematical Statistics, 2004, 342–347. DOI: 10.1214/lnms/1196285403. [Online]. Available: <http://dx.doi.org/10.1214/lnms/1196285403> (cit. on p. 3).
- [19] W. K. Hastings, “Monte carlo sampling methods using markov chains and their applications,” *Biometrika*, vol. 57, no. 1, 97–109, Apr. 1970, ISSN: 0006-3444. DOI: 10.1093/biomet/57.1.97. [Online]. Available: <http://dx.doi.org/10.1093/biomet/57.1.97> (cit. on p. 3).
- [20] S. R. Atance, J. V. Diez, O. Engkvist, S. Olsson and R. Mercado, “De novo drug design using reinforcement learning with graph-based deep generative models,” *Journal of Chemical Information and Modeling*, vol. 62, no. 20, 4863–4872, Oct. 2022, ISSN: 1549-960X. DOI: 10.1021/acs.jcim.2c00838. [Online]. Available: <http://dx.doi.org/10.1021/acs.jcim.2c00838> (cit. on p. 3).
- [21] J. Viguera Diez, S. Romeo Atance, O. Engkvist and S. Olsson, “Generation of conformational ensembles of small molecules via surrogate model-assisted molecular dynamics,” Nov. 2023. DOI: 10.26434/chemrxiv-2023-sx61w. [Online]. Available: <http://dx.doi.org/10.26434/chemrxiv-2023-sx61w> (cit. on p. 4).
- [22] J. Hughes, S Rees, S. Kalindjian and K. Philpott, “Principles of early drug discovery,” *British Journal of Pharmacology*, vol. 162, no. 6, 1239–1249, Feb. 2011, ISSN: 1476-5381. DOI: 10.1111/j.1476-5381.2010.01127.x. [Online]. Available: <http://dx.doi.org/10.1111/j.1476-5381.2010.01127.x> (cit. on p. 5).
- [23] Y. Tabana, D. Babu, R. Fahlman, A. G. Siraki and K. Barakat, “Target identification of small molecules: An overview of the current applications in drug discovery,” *BMC Biotechnology*, vol. 23, no. 1, Oct. 2023, ISSN: 1472-6750. DOI: 10.1186/s12896-023-00815-4. [Online]. Available: <http://dx.doi.org/10.1186/s12896-023-00815-4> (cit. on p. 5).
- [24] V. Pattan, R. Kashyap, V. Bansal, N. Candula, T. Koritala and S. Surani, “Genomics in medicine: A new era in medicine,” *World Journal of Methodology*, vol. 11, no. 5, 231–242, Sep. 2021, ISSN: 2222-0682. DOI: 10.5662/wjm.v11.i5.231. [Online]. Available: <http://dx.doi.org/10.5662/wjm.v11.i5.231> (cit. on p. 5).
- [25] S. Al-Amrani, Z. Al-Jabri, A. Al-Zaabi, J. Alshekaili and M. Al-Khabori, “Proteomics: Concepts and applications in human medicine,” *World Journal of Biological Chemistry*, vol. 12, no. 5, 57–69, Sep. 2021, ISSN: 1949-8454. DOI: 10.4331/wjbc.v12.i5.57. [Online]. Available: <http://dx.doi.org/10.4331/wjbc.v12.i5.57> (cit. on p. 5).
- [26] P. Wah Tang, P. San Chua, S. Kee Chong *et al.*, “A review of gene knockout strategies for microbial cells,” *Recent Patents on Biotechnology*, vol. 9, no. 3, 176–197, Jun. 2016, ISSN: 1872-2083. DOI: 10.2174/1872208310666160517115047. [Online]. Available: <http://dx.doi.org/10.2174/1872208310666160517115047> (cit. on p. 5).

- [27] V. D. Mouchlis, A. Afantitis, A. Serra *et al.*, “Advances in de novo drug design: From conventional to machine learning methods,” *International Journal of Molecular Sciences*, vol. 22, no. 4, p. 1676, Feb. 2021, ISSN: 1422-0067. DOI: 10.3390/ijms22041676. [Online]. Available: <http://dx.doi.org/10.3390/ijms22041676> (cit. on p. 5).
- [28] U. Muhammad, A. Uzairu and D. Ebuka Arthur, “Review on: Quantitative structure activity relationship (qsar) modeling,” *Journal of Analytical & Pharmaceutical Research*, vol. 7, no. 2, Apr. 2018, ISSN: 2473-0831. DOI: 10.15406/japlr.2018.07.00232. [Online]. Available: <http://dx.doi.org/10.15406/japlr.2018.07.00232> (cit. on p. 5).
- [29] P. Gramatica, “Principles of qsar models validation: Internal and external,” *QSAR & Combinatorial Science*, vol. 26, no. 5, 694–701, May 2007, ISSN: 1611-0218. DOI: 10.1002/qsar.200610151. [Online]. Available: <http://dx.doi.org/10.1002/qsar.200610151> (cit. on p. 5).
- [30] D. Weininger, “Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules,” *Journal of Chemical Information and Computer Sciences*, vol. 28, no. 1, 31–36, Feb. 1988, ISSN: 1520-5142. DOI: 10.1021/ci00057a005. [Online]. Available: <http://dx.doi.org/10.1021/ci00057a005> (cit. on p. 7).
- [31] V. G. Satorras, E. Hoogeboom, F. B. Fuchs, I. Posner and M. Welling, *E(n) equivariant normalizing flows*, 2022. arXiv: 2105.09016 [cs.LG] (cit. on pp. 8–10).
- [32] M. Geiger and T. Smidt, *E3nn: Euclidean neural networks*, 2022. arXiv: 2207.09453 [cs.LG] (cit. on pp. 8, 10).
- [33] K. T. Schütt, P.-J. Kindermans, H. E. Sauceda, S. Chmiela, A. Tkatchenko and K.-R. Müller, *Schnet: A continuous-filter convolutional neural network for modeling quantum interactions*, 2017. arXiv: 1706.08566 [stat.ML] (cit. on p. 8).
- [34] T. S. Cohen and M. Welling, *Group equivariant convolutional networks*, 2016. arXiv: 1602.07576 [cs.LG] (cit. on p. 8).
- [35] S. L. Batzner, A. Musaelian, L. Sun *et al.*, “E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials,” *Nature Communications*, vol. 13, 2021. [Online]. Available: <https://api.semanticscholar.org/CorpusID:231418897> (cit. on p. 8).
- [36] J. Zhou, G. Cui, S. Hu *et al.*, *Graph neural networks: A review of methods and applications*, 2021. arXiv: 1812.08434 [cs.LG] (cit. on p. 8).
- [37] B. Khemani, S. Patil, K. Kotecha and S. Tanwar, “A review of graph neural networks: Concepts, architectures, techniques, challenges, datasets, applications, and future directions,” *Journal of Big Data*, vol. 11, no. 1, Jan. 2024, ISSN: 2196-1115. DOI: 10.1186/s40537-023-00876-4. [Online]. Available: <http://dx.doi.org/10.1186/s40537-023-00876-4> (cit. on p. 8).

- [38] R. Mercado, T. Rastemo, E. Lindelöf *et al.*, “Graph networks for molecular design,” Aug. 2020. DOI: 10.26434/chemrxiv.12843137.v1. [Online]. Available: <http://dx.doi.org/10.26434/chemrxiv.12843137.v1> (cit. on pp. 9, 10, 13).
- [39] J. Guan, W. W. Qian, X. Peng, Y. Su, J. Peng and J. Ma, *3d equivariant diffusion for target-aware molecule generation and affinity prediction*, 2023. arXiv: 2303.03543 [q-bio.BM] (cit. on pp. 9, 10).
- [40] N. Brown, M. Fiscato, M. H. Segler and A. C. Vaucher, “Guacamol: Benchmarking models for de novo molecular design,” *Journal of Chemical Information and Modeling*, vol. 59, no. 3, 1096–1108, Mar. 2019, ISSN: 1549-960X. DOI: 10.1021/acs.jcim.8b00839. [Online]. Available: <http://dx.doi.org/10.1021/acs.jcim.8b00839> (cit. on p. 9).
- [41] O. Engkvist, P.-O. Norrby, N. Selmi *et al.*, “Computational prediction of chemical reactions: Current status and outlook,” *Drug Discovery Today*, vol. 23, no. 6, pp. 1203–1218, 2018, ISSN: 1359-6446. DOI: <https://doi.org/10.1016/j.drudis.2018.02.014>. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1359644617305068> (cit. on p. 9).
- [42] C. W. Coley, W. H. Green and K. F. Jensen, “Machine learning in computer-aided synthesis planning,” *Accounts of Chemical Research*, vol. 51, no. 5, 1281–1289, May 2018, ISSN: 1520-4898. DOI: 10.1021/acs.accounts.8b00087. [Online]. Available: <http://dx.doi.org/10.1021/acs.accounts.8b00087> (cit. on p. 9).
- [43] C. W. Coley, D. A. Thomas, J. A. M. Lummiss *et al.*, “A robotic platform for flow synthesis of organic compounds informed by ai planning,” *Science*, vol. 365, no. 6453, Aug. 2019, ISSN: 1095-9203. DOI: 10.1126/science.aax1566. [Online]. Available: <http://dx.doi.org/10.1126/science.aax1566> (cit. on p. 9).
- [44] A. Tayyebi, A. S. Alshami, Z. Rabiei *et al.*, “Prediction of organic compound aqueous solubility using machine learning: A comparison study of descriptor-based and fingerprints-based models,” *Journal of Cheminformatics*, vol. 15, no. 1, Oct. 2023, ISSN: 1758-2946. DOI: 10.1186/s13321-023-00752-6. [Online]. Available: <http://dx.doi.org/10.1186/s13321-023-00752-6> (cit. on p. 9).
- [45] L. Keutzer, H. You, A. Farnoud *et al.*, “Machine learning and pharmacometrics for prediction of pharmacokinetic data: Differences, similarities and challenges illustrated with rifampicin,” *Pharmaceutics*, vol. 14, no. 8, p. 1530, Jul. 2022, ISSN: 1999-4923. DOI: 10.3390/pharmaceutics14081530. [Online]. Available: <http://dx.doi.org/10.3390/pharmaceutics14081530> (cit. on p. 9).
- [46] J. Köhler, A. Krämer and F. Noé, *Smooth normalizing flows*, 2021. arXiv: 2110.00351 [stat.ML] (cit. on pp. 9, 11).
- [47] M. Dibak, L. Klein, A. Krämer and F. Noé, *Temperature steerable flows and boltzmann generators*, 2022. arXiv: 2108.01590 [cond-mat.stat-mech] (cit. on p. 9).

- [48] H. Wu, J. Köhler and F. Noé, *Stochastic normalizing flows*, 2020. arXiv: 2002.06707 [stat.ML] (cit. on pp. 9, 11).
- [49] E. Mansimov, O. Mahmood, S. Kang and K. Cho, “Molecular geometry prediction using a deep generative graph neural network,” *Scientific Reports*, vol. 9, no. 1, Dec. 2019, ISSN: 2045-2322. DOI: 10.1038/s41598-019-56773-5. [Online]. Available: <http://dx.doi.org/10.1038/s41598-019-56773-5> (cit. on p. 9).
- [50] O.-E. Ganea, L. Pattanaik, C. W. Coley *et al.*, *Geomol: Torsional geometric generation of molecular 3d conformer ensembles*, 2021. arXiv: 2106.07802 [physics.chem-ph] (cit. on p. 9).
- [51] M. Xu, L. Yu, Y. Song, C. Shi, S. Ermon and J. Tang, *Geodiff: A geometric diffusion model for molecular conformation generation*, 2022. arXiv: 2203.02923 [cs.LG] (cit. on p. 9).
- [52] B. Jing, G. Corso, J. Chang, R. Barzilay and T. Jaakkola, *Torsional diffusion for molecular conformer generation*, 2023. arXiv: 2206.01729 [physics.chem-ph] (cit. on pp. 9, 15).
- [53] J. Jiménez-Luna, F. Grisoni and G. Schneider, “Drug discovery with explainable artificial intelligence,” *Nature Machine Intelligence*, vol. 2, no. 10, 573–584, Oct. 2020, ISSN: 2522-5839. DOI: 10.1038/s42256-020-00236-4. [Online]. Available: <http://dx.doi.org/10.1038/s42256-020-00236-4> (cit. on p. 9).
- [54] P. Linardatos, V. Papastefanopoulos and S. Kotsiantis, “Explainable ai: A review of machine learning interpretability methods,” *Entropy*, vol. 23, no. 1, 2021, ISSN: 1099-4300. DOI: 10.3390/e23010018. [Online]. Available: <https://www.mdpi.com/1099-4300/23/1/18> (cit. on p. 9).
- [55] S. Viet Johansson, H. Gummesson Svensson, E. Bjerrum *et al.*, “Using active learning to develop machine learning models for reaction yield prediction,” *Molecular Informatics*, vol. 41, no. 12, Jul. 2022, ISSN: 1868-1751. DOI: 10.1002/minf.202200043. [Online]. Available: <http://dx.doi.org/10.1002/minf.202200043> (cit. on p. 9).
- [56] H. G. Svensson, C. Tyrchan, O. Engkvist and M. H. Chehreghani, *Utilizing reinforcement learning for de novo drug design*, 2023. arXiv: 2303.17615 [q-bio.BM] (cit. on p. 9).
- [57] J. Yosinski, J. Clune, Y. Bengio and H. Lipson, *How transferable are features in deep neural networks?* 2014. arXiv: 1411.1792 [cs.LG] (cit. on p. 10).
- [58] M. E. Peters, S. Ruder and N. A. Smith, *To tune or not to tune? adapting pretrained representations to diverse tasks*, 2019. arXiv: 1903.05987 [cs.CL] (cit. on p. 10).
- [59] M. Moret, L. Friedrich, F. Grisoni, D. Merk and G. Schneider, “Generative molecular design in low data regimes,” *Nature Machine Intelligence*, vol. 2, no. 3, 171–180, Mar. 2020, ISSN: 2522-5839. DOI: 10.1038/s42256-020-0160-y. [Online]. Available: <http://dx.doi.org/10.1038/s42256-020-0160-y> (cit. on p. 10).

- [60] V. D. Mouchlis, A. Afantitis, A. Serra *et al.*, “Advances in de novo drug design: From conventional to machine learning methods,” *International Journal of Molecular Sciences*, vol. 22, no. 4, p. 1676, Feb. 2021, ISSN: 1422-0067. DOI: 10.3390/ijms22041676. [Online]. Available: <http://dx.doi.org/10.3390/ijms22041676> (cit. on p. 10).
- [61] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, no. 8, 1735–1780, Nov. 1997, ISSN: 1530-888X. DOI: 10.1162/neco.1997.9.8.1735. [Online]. Available: <http://dx.doi.org/10.1162/neco.1997.9.8.1735> (cit. on p. 10).
- [62] J. Lim, S. Ryu, J. W. Kim and W. Y. Kim, “Molecular generative model based on conditional variational autoencoder for de novo molecular design,” *Journal of Cheminformatics*, vol. 10, no. 1, Jul. 2018, ISSN: 1758-2946. DOI: 10.1186/s13321-018-0286-7. [Online]. Available: <http://dx.doi.org/10.1186/s13321-018-0286-7> (cit. on p. 10).
- [63] D. P. Kingma and M. Welling, *Auto-encoding variational bayes*, 2022. arXiv: 1312.6114 [stat.ML] (cit. on p. 10).
- [64] L. Maziarka, A. Pocha, J. Kaczmarczyk, K. Rataj, T. Danel and M. Warchol, “Mol-cycleGAN: A generative model for molecular optimization,” *Journal of Cheminformatics*, vol. 12, no. 1, Jan. 2020, ISSN: 1758-2946. DOI: 10.1186/s13321-019-0404-1. [Online]. Available: <http://dx.doi.org/10.1186/s13321-019-0404-1> (cit. on p. 10).
- [65] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza *et al.*, *Generative adversarial networks*, 2014. arXiv: 1406.2661 [stat.ML] (cit. on p. 10).
- [66] J. Zhou, G. Cui, S. Hu *et al.*, *Graph neural networks: A review of methods and applications*, 2021. arXiv: 1812.08434 [cs.LG] (cit. on p. 10).
- [67] J. You, B. Liu, R. Ying, V. Pande and J. Leskovec, *Graph convolutional policy network for goal-directed molecular graph generation*, 2019. arXiv: 1806.02473 [cs.LG] (cit. on p. 10).
- [68] W. Feng, L. Wang, Z. Lin *et al.*, “Generation of 3d molecules in pockets via a language model,” *Nature Machine Intelligence*, vol. 6, no. 1, 62–73, Jan. 2024, ISSN: 2522-5839. DOI: 10.1038/s42256-023-00775-6. [Online]. Available: <http://dx.doi.org/10.1038/s42256-023-00775-6> (cit. on p. 10).
- [69] J. Ho, A. Jain and P. Abbeel, *Denoising diffusion probabilistic models*, 2020. arXiv: 2006.11239 [cs.LG] (cit. on p. 10).
- [70] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon and B. Poole, *Score-based generative modeling through stochastic differential equations*, 2021. arXiv: 2011.13456 [cs.LG] (cit. on p. 10).
- [71] Y. Song, C. Durkan, I. Murray and S. Ermon, *Maximum likelihood training of score-based diffusion models*, 2021. arXiv: 2101.09258 [stat.ML] (cit. on p. 10).
- [72] G. Papamakarios, E. Nalisnick, D. J. Rezende, S. Mohamed and B. Lakshminarayanan, *Normalizing flows for probabilistic modeling and inference*, 2021. arXiv: 1912.02762 [stat.ML] (cit. on p. 10).

- [73] Z. Wu, B. Ramsundar, E. Feinberg *et al.*, “Moleculenet: A benchmark for molecular machine learning,” *Chemical Science*, vol. 9, no. 2, 513–530, 2018, ISSN: 2041-6539. DOI: 10.1039/c7sc02664a. [Online]. Available: <http://dx.doi.org/10.1039/C7SC02664A> (cit. on p. 11).