

# Structure-Based Molecular Optimisation Through Deep Generative Networks: AI-Enhanced Computational Strategies for Rational Drug Design

*Dr. Ying Liu, Associate Professor of Computer Science, Nanyang Technological University (NTU), Singapore*

---

---

## 1. Introduction

Chemical biology can't escape from the marriage between computational science and high-throughput experimentation, either intentionally or unintentionally. The large amounts, variability, higher dimensions, and flexibility of data are supported and applied through model-driven AI techniques, producing clinical interest modules represented as models implicated in pathological situations or patient responses, such as bioactivity, imaging, and mathematical models, blending strong learning from target and clinical assays with high-throughput contextual data.

HRAN, an emerging technology revolution, is an AI-powered learning approach to computational drug design that has a profound impact on the entire drug discovery and development journey, improving the evolving complex scientific tasks that still need to be carried out at each step.

Over the past few decades, it has generally been difficult to design an effective chemical on the intended target due to the trade-offs between diversity and synthetic accessibility, the high cost dedicated to redundant and irrelevant chemical modifications that have already been pursued, low-throughput virtual screening and docking methods that generate few appealing hits, and the failure to handle increasingly large and fuzzy chemical data while still maintaining well-calibrated and generalizable learning models. What differentiates the accessible pool of active molecules is typically sensitive changes in the processing and experimental methodologies that have become routine in the medicinal chemistry group acting as an artificial chemist, as opposed to the robust predictive models validated by many groups. Many large pharmaceutical corporations

are conducting similar studies. This is an interesting area of development in the pharmaceutical industry due to the large amount of literature and the quantitative structure–activity relationship cheminformatics. The combination of sophisticated AI systems with design ideas could be instrumental in raising the temporary medicines and energy levels required for effective development. It is our expectation that this chapter will promote collaboration between chemists and researchers by showcasing new strategies like machine learning, crowdsourced data visualizations, and automated model development and testing using big data technologies, which may provide a solution for interdisciplinary data and knowledge discovery to increase the positive effects of drug development.

### **1.1. Background and Significance**

The process of identifying, lead optimization, preclinical testing, clinical development, regulatory approval, and market release takes an average of 14 years, and the cost is typically more than ninety million for each therapeutic approved. Not only is the number of compounds approved per year small, but the full spectrum of drug development takes an average of 20 years. The computational methods based on physicochemical properties were first introduced to drug design in the 1970s, and their various developments can be divided into two broad methods: molecular mechanics-based drug design and structure-based drug design. Lead optimization covers high-throughput computational chemical reactions and excludes compound data improvement in terms of absorption, distribution, metabolism, excretion, and toxicity profiling, which currently is the cornerstone of drug development.

Artificial intelligence has achieved success in various domains in recent years. In particular, computer vision and natural language processing perform at human level. Integrating AI in drug discovery has become the hottest issue in the academic and industrial community since it profoundly improves the processes in the early drug development stage. A number of advances and significant achievements have been made in a small number of reports researching various aspects of drug discovery, such as compound design, de novo drug generation, structure-activity relationship prediction, and auto-generated lead optimization. Current AI can build on the attractive pathway for addressing various problems in traditional drug discovery. Chemical space may conduce to billions or trillions of candidates on various fragments and molecular

structures that can be designed by AI. Rapid growth of medicinally relevant data has been collected to generate, develop chemoinformatics, and provide alerts associated with drug safety. The boom of high-throughput inhibitor screens and X-ray crystallography has built a narrow ligand and target-based computational approach in structure-based drug design. This experiment provides reasonably appropriate computation for predicting structural changes in the target caused by ligand binding. Several case studies have seen the integration of computational experimental methods, and it has led to the discovery of very potent drug candidates. The AI-concealed problems will consequently affect the difficulties in the early stages of drug discovery, hence reducing the long-lasting and costly process. In the long run, public benefits include completely novel therapeutics and addressing patient suffering and societal expenses for chronic disorders like psychiatric illnesses, diabetes, and cancer.

## **1.2. Research Objectives**

The main research objective is to study whether and under what conditions the AI methodology can be applied to improve the drug development process. Specifically, we seek to provide some answers to the following, more detailed questions:

- a) Can machine learning, combined with computational methodologies, help in obtaining better predictions of the interactions of drug compounds with specific biochemical targets? Here, we will be interested in the impacts of extending the applicability of the corresponding methods;
- b) Can we improve the selection of a 'lead' compound from a number of diverse candidates? Lead compounds are the first important step in drug discovery: they are potential 'starter' new drugs for testing and analysis. Once more, the intention is to consider if machine learning can accelerate current computational approaches, in this case by predicting toxicity issues for the final drug application, as early in the part of the drug development process as possible;
- c) Can we integrate into the drug design equation more qualitative AI results such as drug-likeness, or penalize agents that need a more lengthy development cycle because of their underlying physicochemical properties? Thus, a second main objective will be to explore the potential for integrating machine learning methods within the computational drug design process;

d) We will also rigorously test and evaluate any potential improvements from these initial investigations. Although qualitative results are difficult to evaluate quantitatively, careful testing will be carried out and actual improvements observed and documented. The goal of many scientific endeavors is to evaluate a hypothesis, or hypotheses, about a system being studied. In computational drug design, which is typically in data-poor and information-limited environments, how to quantify and assess the likelihood of work being 'successful' is less straightforward, and so a main final intrinsic objective will be to define critical quality indicators or benchmark criteria to assess the overall impact of machine learning approaches to drug design, and identify stronger qualitative results or reasonably quantify the enhancement of key computational tools to measure the nurture of new technologies.

Finally, our research may also indirectly propose answers to the following intriguing questions: what, overall, is the relative cost of using machine learning tools in conjunction with a wide variety of computational drug design approaches likely to bring about? This is a critically important question to healthcare providers whose primary task is to deliver effective care affordably. More precise, targeted, and efficient drugs, with reduced numbers of associated side effects, mean a reduction in hospital budgets dedicated to drug-induced diseases and processes - whether this is on equipment, personnel, or staff retraining and readaptation.

## **2. Fundamentals of Computational Drug Design**

Computational drug design has become a critical approach for identifying and designing biologically active compounds, especially with the recent surge in demand for pharmaceuticals. At the core of this field of research is the integration of computers into most, if not all, processes, giving rise to the moniker "in silico" methods. The high cost and time needed for experimental procedures have made in silico methods an attractive choice for most pharmaceutical R&D ventures. For computational drug design to be effective, it is critical for any biological system to be simulated accurately, necessitating knowledge of biomolecular structures as well as ligand-receptor interactions. A variety of software tools and algorithms are essential to the structure-assisted drug design pipeline. State-of-the-art strategies in structure-assisted drug design include virtual screening, which ranges from docking-based methods to ligand-based methods such as quantitative structure-activity relationships and pharmacophore modeling.

Computational methods are complementary to the work done through extensive experimentation, and together they yield excellent results in drug discovery. High-throughput experimentation, synthesis, and screening of hundreds of substrates at a time have become possible due to fast-paced computational methods. High-throughput screening of ligand molecules is more effective and requires less time and fewer expenses. The structure of macromolecules is the main focus of material scientists and biologists who use computational methods to investigate the functional role of these molecules and to study their unique properties. Computational methods are used in preclinical trials to test the performance of potential drug candidates. The physical and chemical properties of potential drug candidates are evaluated based on sophisticated computational tools before they are advanced to the trial phases.

### **2.1. In Silico Drug Screening**

In silico drug screening, also known as in silico virtual screening, has made remarkable progress in recent years, particularly as computational tools and algorithms have continued to mature and expand. These techniques can be used to focus the attention of researchers on drug-like compounds from very large libraries using a wide array of computational models. Using only a candidate pharmacophore, an array of in silico tools can rapidly amass large quantities of second-generation substance evaluations thanks to the availability of molecular structural data stored in databases. While derived from a much simpler in vitro strategy for fruiting couples, in vitro pharmacophore modeling has not yet replaced VS strategies based on ligand–receptor interaction data.

Consequently, in silico techniques have the effectiveness to streamline the candidate schedule to practical constants while lowering the cost of achieving a specific compound, all while facing a significant rise in pharmaceutical research and development expenditures. One of the biggest hurdles in computer-aided molecular design is the elevation of pure data into meaningful screening results. To overcome these issues, in silico programs need several refined preprocessing steps to cope with the multi-information layer nature of these large databases. However, refining strategies to improve the quality of these datasets is a significant challenge when working with datasets of any extent. Those additional preprocessing or cleaning steps are only necessary if they deal with particularly high disorder. These requirements limit data quality and prevent noise-induced random educated guesses and ensure structures from

having a positive spectral signature. Computation-based drug screening complements and assists rather than replaces advanced biological testing. Electronic-based drug screening is an integral component of what has become an increasingly iterative approach to drug discovery.

## **2.2. Molecular Docking Techniques**

Molecular docking is a computational technique for studying and simulating molecular interactions, including protein-ligand interactions, and remains a powerful tool in rational computational drug design. It aims to predict the ligand-receptor binding geometry and affinity or poses that are thermodynamically preferable for the complex. It adopts various computational algorithms and theories that work together to produce physiologically relevant binding affinities between ligands and protein receptors. From a methodological perspective, a number of steps can be identified that are commonly used in docking procedures. The overall procedure, therefore, consists of sampling binding configurations, optimizing structure conformations, and evaluating pose binding affinities.

Given the central role of docking techniques in drug design, a large number of software tools have been developed over the years. Docking simulations are performed with some of them focusing on the optimization of biological drug-receptor processes at the atomic level, i.e., by paying particular attention to the atomic interactions that govern the mechanism on a biophysical level. Most of these tools follow different methodologies and thus have both limitations and particular strengths in drug discovery. In general, the great majority of software tools providing molecular docking facilities are not specifically designed for distinguishing between types of molecular interactions with more than an atomistic level description, while they do not even take into account the contribution of the solvent during the process. In a simulation which allows, typically by the Langevin thermostat, the protein was embedded to simulate the interaction with the surrounding water.

Several approaches have been used to perform the actual docking, including predominantly Lennard-Jones and Coulomb interactions, but also van der Waals interactions, and variations with more stages of free energy can achieve a better performance. In a thorough assessment and validation of some of the most common root mean square deviation results, more than 25% immobility in the docking protein is in

the range of user-specified parameters at best, moving typically applied restraints to the initial crystallographic rigid motion of protein hydrogens. In this approach, extensive manual validation of the protein-ligand complexes was also performed, with the observation that a great majority of the compounds had multiple poses within similar binding energy scores to the deposited one within crystallography. This underlines the potential in using molecular dynamic simulation to further explore the degree of specificity of an interaction between a target and a ligand. The necessary flexibility allows the ligand to adapt to the environment from the related proteins of the binding modes needed to search. The presence of binding sites in many physiologically active proteins is flexible in response to the presence of light. Frequently, only a part of the receptor contributes to the initial ligand-receptor complex. It is often the first site to incorporate the binding energy for further elucidation.

### **3. Machine Learning in Drug Discovery**

3.1 Introduction to Machine Learning in Drug Discovery Machine learning is a type of artificial intelligence where computer algorithms learn from and make predictions, recommendations, and decisions from data. Machine learning is particularly beneficial in high-dimensional problems where traditional statistical methods struggle, as is typically the case in drug discovery. The application of machine learning in drug discovery is varied and usually serves as a supportive tool for the prediction of the targets of small molecules, off-target binding of candidate therapeutics, drug toxicity, and general drug safety prediction. In silico absorption, distribution, metabolism, excretion, and toxicology analyses of drug-like compounds are also common applications of machine learning in drug discovery.

3.2 Machine Learning in Drug Design Machine learning methods can offer faster and cheaper alternatives to traditional computationally expensive molecular dynamics, molecular docking, and pharmacophore modeling approaches, which have played a crucial role in drug design for decades. By employing a wide range of machine learning approaches as the predictive model, in silico calculations of the properties and interactions of molecules with biological targets can be expedited with moderate performance levels. In this context, structure-based virtual screening for drug discovery and development can be made much more efficient by using machine learning prediction models, which can review millions of molecules in vastly shorter periods.

Machine learning techniques commonly employed in drug discovery to handle and analyze molecular, pharmacological, and clinical data include supervised learning, used to ascertain correlations between a set of input features and the prediction target. Unsupervised learning models identify patterns or categorize data points based on the characteristics of the input features and can also be employed to provide insight into the relative importance of the individual input features in a supervised learning model. Semi-supervised and reinforcement learning models have been less frequently used within the context of drug discovery requirements due to the relative lack of data points and computationally intensive learning processes, respectively. Utilization of machine learning models in drug discovery is also hindered by the modest size of available datasets, their quality, as well as the intra- and inter-lab variability in protocols utilized for generating this data.

### **3.1. Supervised Learning Algorithms**

Thus far, the work using AI in efficient drug discovery and development has strictly been classified under supervised learning approaches. The reason for this is simple: the only requisite input for supervised learning methods is the available experimental data for the molecules to which the anticipated drug groomers will be exposed. That "training" molecule set is one case of the more general approach of using computational methods to identify the distinguishing molecular features of different chemical spaces, i.e., cheminformatics. Ideally, accurate models may be created with a small, highly diverse set of molecules to train the supervisor, and computing time grows linearly with the number of molecules in the chemical space.

This is, of course, only one of the optimistically placed but unrealized potentials of cheminformatics. More limited success has been achieved by predicting chemical biomimetic behavior with supervised computational methods. These include high-throughput screening, virtual screening, and ligand-receptor molecular docking (key determinants of the success of these methods are the composition of the training molecule sets used, i.e., their similarity of training sets to the respective test set). In the world of bioprospecting, prediction of specific chemical properties emergent in a phylogenetic clade could be used to guide the discovery of novel therapeutic compounds from within poorly studied chemical taxa. This mentions the key aspects of a successful supervised chemical training program: validated empirical relationships

between the structural attributes of the training molecules and one or more attributes of the test molecules.

### **3.2. Unsupervised Learning Techniques**

Unsupervised learning (UL) aims to find hidden structures compressed in input data, such as searching for interesting clusters, finding lower-dimensional representations like principal component analysis, or converting input data to a form that distinguishes relevant features of the given data, such as manifold learning. As opposed to supervised learning methods, the supervised framework maps the input to scalar output or class label pairs. The comparative analysis of supervised and unsupervised learning is summarized in a table. In supervised learning, the output labels are assumed to be provided to supervise the framework.

The specific task is to automatically derive the class label from a set of features to which it has been mapped during the training of the model. In unsupervised learning, the task is to find structure in the feature representation itself. Examples of unsupervised learning include principal component analysis, which maps the given set of input data into the most principal directions, i.e., eigenvectors, by maximizing input variance; clustering algorithms such as K-means clustering, which groups similar input samples based on distance metrics; and stochastic neighbor embedding, which maps the input pairwise similarities to the low-dimensional space for visualization.

### **4. Integration of AI and Computational Drug Design**

Docking computational drug design methods are limited by the exponential growth of the number of virtual molecules, creating a necessity for proper focusing of the search. In this sphere, AI is at the heart of enabling identification of active compounds through the production of compound libraries using generative models, as well as supporting quick heavy similarity measurements. Furthermore, optimization of compounds themselves can be done by algorithms searching for highly affine chemical compounds based on highly dependent deep learning drug response relationships. AI can effectively rethink molecular simulations, helping to realize novel drug candidates within the first year. Additionally, research shows the possibility to link pharmacodynamics and pharmacokinetics, and association with causation. In some computational methods, AI algorithms can operate on various levels. Interestingly, an interaction between

traditional computational methods and AI-enhanced molecular simulations can provide higher results.

Integration of AI is relatively hand-in-hand with computational drug design exploitation; however, it is not without complications. Using neural networks requires fluency in scripting, non-AI-enhanced methodology, pre-engineering the inputs for architecture predictions, and making specific installments in the hardware. Therefore, additional computing power, a larger research staff training period, and additional computational drug design methodologies are created. Just as potential risks may be seen as a barrier to the effective integration of AI drug design, they offer strong opportunities by bridging the original computational drug design results faster and more accurately. In addition, a proper combination of AI and deep learning will create immense amounts of biases, speeding up the process of searching for new drug candidates. Such a comprehensive process of exploring and enriching the practical results of the drug development process justifies vast collaborations of research team members from the medicine, computer science, biology, and physics-biochemistry areas.

#### **4.1. Challenges and Opportunities**

Despite the fast pace in the field of AI/ML, there are several challenges and opportunities to be faced when integrating AI/ML into computational drug design. On the one hand, big data on biological activities is not necessarily reliable, and there is the issue of one museum/many drugs; that is, all drugs potentially target multiple proteins given genetic interference in proteomes. Since graduality in the association between inputs and responses underlies the robustness and adaptability of life, the use of oversimplistic machine learning algorithms to bypass this complexity hides systemic failure under pseudo-predictive power. This is particularly a concern when black-box models are used that can lead to the situation where reality-check capacity is lost. Another challenge comes from the issue of 'complex systems' where a genome-encoded protein is only one bridge in the pathophysiological association networks that the pharma-driven fantasy of 'on-target' purity neglected or oversimplified. Systems biology is instrumental to move from 'one drug/one target' toward 'network pharmacology', and 'prediction' might only be improved through partnerships with experimental scientists. However, AI might offer an improved use of big data and reduce opinions in silico on what to expect or to screen next for experimental validation.

In undoubtedly good fashion, AI can, where available, analyze big databases and reveal otherwise unpredicted correlations between variables that reflect the structural and functional regulation of some bioactivity. However, the ability of AI to extrapolate from such datasets is questionable given it is highly probable that it would not have learned the difference between 'systemic integration' and 'artificial separations' performed in many experimental setups. AI algorithms are also prone to repetitive and/or similar errors accumulated through subsequent process steps, given they optimize a certain 'objective function', early informed by how to weigh/optimize inputs but still unable to learn what needs to be a systemic—perhaps most things—beyond assisted processes, etc. Consequently, drugs may be shown to pass 'validation' steps despite not being validated in terms of enhancing health, curing, preventing, and prolonging life. Ongoing research and development have to clarify and address these challenges and opportunities. This might well turn other issues of difficulty into potentially high selective advantages. But as trendy AI melds with biology, new challenges and opportunities arise too for computational drug design.

## **5. Case Studies and Applications**

There are a number of case studies that present many examples showing how AI in drug discovery will revolutionize different aspects of the current pharmaceutical process. All are fueled by successful case studies briefly reviewed in each of the subsequent examples of applications exemplified in real-world AI augmentation of the drug discovery pipeline. 1. Improved Target Identification 2. Accelerated Compound Optimization 3. Improved Safety Assessment One of the lessons learned from these case studies is that projects of this magnitude need to be conducted by a multidisciplinary team where computational chemists lead a group of experts representing AI, informatics, and pharmacoeconomics. Secondly, the significance of simply publishing or being part of a publication in a high-ranking impact journal is not a goal in this effort. It is the capability of the work to fuel company innovation. This emphasizes the area of AI-augmented precision medicine in the field of cardiovascular diseases and CNS. Along with eye coordination, we hope to ensure the safety pharmacological profile of new leads collected in the initial steps of our AI-dedicated efforts in the discovery of new compounds. This showcases the actual impact of AI augmentation on the traditional drug discovery process. The majority of areas of drug discovery are prone to AI

harvesting, thus causing an actual paradigm shift in our lives, pointing to the ways we work. Especially, we hope that the sales of reagents for HTS will drop by 50% by 2025.

### **5.1. Successful Implementations**

Machine learning (ML) and more generally AI have a significant impact on drug molecule discovery. Progressive strategies that drive bottom-line value in AI-augmented drug discovery are gradually taking shape. Demonstrative case studies are presented describing successful applications of in silico methods for the computational design of therapeutic compounds. These implementations illustrate in very practical terms what impacts the AI augmentation of drug development may offer. Notably, the three AI case studies presented depict the earliest associated reduction of lead generation time, up to the discovery of drug candidates reaching clinical trials. The lesson that can be learned from these illustrative case studies must be put cautiously. AI contributions carry with them the authority of group visibility: the expert-driven choices and decisions, parallel data and experimental and informatics effort, subsequent system design, and skillful hardcore engineering aspects.

Thanks to breakthroughs in AI technologies, it is now commercial to vet these scientifically creative insights in robust ML systems at multiple stages of a medicinal chemistry program. An ML-guided approach can supersede pathway brainstorming with data, to highlight beneficial synthesis sequences, knock-out approaches, pruning of blind alleys, and re-prioritize long-term decision making. The impact of such a transformative molecule design methodology in significantly enhancing future clinical success rates can be expected to far outstrip the value of modest improvements in AI-optimized strategies within lead optimization. At the core of applying AI in computer-aided drug design are cutting-edge methods both in solo AI and molecular design. These have already been implemented in cases of disease areas and types such as oncology, dermatology, or antibiotic-resistant bacteria. Such approaches build on a robust data set and expert-driven rules that inform ML and AI designs, which are iteratively evaluated and tested. The most promising top hits are subsequently validated experimentally for robustness and performance. The AI aims to be compatible with state-of-the-art strategies for disease treatment. This ensures that the proposed designs have an immediate benefit for target populations and orphan diseases and provide translational potential with strong normative regulatory foundations.

## 5.2. Future Directions

The rising body of work in the AI-driven drug discovery space is extremely promising, but also serves to further validate the dynamism and adaptability of the field. What follows are a number of future directions that we anticipate taking this area of research over the next few years—one that is both as invigorating and unsettling as it is speculative.

The application of AI to drug design is currently focused on major therapeutic areas: for which generally clear points of intervention are known. There is a forthcoming and unique opportunity to apply AI in those smaller yet validated areas of investigation—like analgesics and corticosteroids for major depression that can be directionally based on preliminary data. Further, it is dawning upon researchers that the largest remaining questions to be answered in molecular and cellular sciences are to be solved computationally. The data per publication are significantly decreasing—a sign of our alignment with the limits of conventional goals—while at the same time, the resolution of the human body, for example to access physiologically therapeutic concentrations of drug in relatively healthy tissue, is surpassing capabilities: meaning individual/private field robotics are accelerating at an exponential rate. The building blocks are therefore in place for personalized in silico modeling of an individual, opening an area of drug design we can hardly fathom at this time.

The first regulatory approved systems in this area will have proprietary privilege for a period. We have already seen some entree to the utilization and co-utilization of quantum computing in conjunction with traditional computing, and we expect to see advanced data analytic platforms provide similar “leaps” in the near future to take drug discovery into new and innovative applications of AI models. Pharmaceutical companies have been slow to adapt an agile informatic approach, however the expectation is that new AI product successes are going to attack and offset traditional capabilities consecutively with each new development. In short, there is no parameter in this aspect of the drug discovery game that will not change. The combination of AI adopting these strategies and increasing speed of capability advancement while dropping in price point to implementation—and convergence upon our endgame that detailed our choice to create a foundation of Coins or Currency by a drug-based ecommerce platform—is perilously close.

## 6. Conclusion

New technologies boost us into a new era of various fields, which has been more efficient and capable. One of the most advanced fields is drug design, where one of the most successful technologies is AI. The AI-based techniques have been developed using the technologies used for chemical- or biology-based studies. While the other fields have mainly focused on a specific field containing AI, those studies used the combination of either sophisticated computational methods and data or only data from the computational studies. Considering the importance of the effects of the behaviour of AI, most studies combined machine learning and the computational approach. This approach provides successful and significant outcomes. Furthermore, the study areas of AI-based computational methods combined drug design are spread across the candidate identification to ethical considerations. The current available drugs are also retried and repurposed to design better affects on patients. Therefore, the integration uses automatic approaches, especially the reinforcement self-learning technique. These techniques are able to develop a wide range of potential drug candidates using the approved drugs and newly developed compounds. Several machine learning and computational techniques provide us with potential tools for ADE pattern prediction, which is constructed by the given compound from the property- and target-based studies. This summary introduces the machine learning and computational drug candidates are used from the wet-lab study to in silico calculation. To summarise, the integrated use of AI in drug design has allowed a tremendous reduction in time and cost for new drug development. Despite these strides in the field, there are yet to be developed drug candidates from AI and thus they have never been tested clinically. So many challenges and opportunities have yet to be addressed in the AI-based method applications. Machine learning, especially the deep learning from molecular structure, also needs to go deeper, focusing on the biological and informative aspects of drug design to make patient care effective and safe.