# COMMENTARY

# Redefining drug discovery is deep generative molecular design

Matt Coban, Christopher Hopkins

## ABSTRACT

Deep generative models and Artificial Intelligence (AI) have recently made strides that have demonstrated their value in the medical field, particularly in the drug discovery and development process. The developer and user must decide which protocols to take into account, which elements to carefully examine, and how deep generative models may combine the necessary disciplines in order to properly deploy AI. This study provides an updated and user-friendly reference for the large computational drug discovery and development community by summarizing traditional and recently emerging AI methodologies. From various angles, we introduce deep generative models and discuss the theoretical underpinnings of describing chemical and biological structures as well as their practical applications. We go over the data issues and technical difficulties and highlight the multimodal deep generative models potential for speeding up drug discovery.

## INTRODUCTION

According to a recent report, pharmaceutical corporations spent $2.6 billion in 2015—an increase from $802 million in 2003—on the development of new, US Food and Drug Administration-approved medications. Although there are higher direct expenditures associated with clinical trials, the capitalized costs of the two stages are essentially comparable because the preclinical investment occurs earlier. The requirements and urgencies are captured by recent developments in computational sciences and technology, which also offer a number of potentially fruitful strategies. Developers can choose the best Artificial Intelligence (AI) among these to tackle the issue at hand, particularly deep generative models, suitable protocol, and factors. Together, they illustrate pathways that unite pharmacology, computer science, biology, chemistry, and strategies for treating diseases. AI for drug development has made significant strides thanks to the quick increase in processing power, data volume, and sophisticated algorithms, particularly in the use of deep generative models. The models have shown great promise in transforming small-molecule and macromolecule design, optimization, and synthesis. Deep generative model applications have already produced new partially optimized candidate leads, sometimes in a fraction of the time needed by traditional sequential procedures. Deep generative modeling has the potential to speed up the Research and Development (R&D) process if used widely. With the use of data structures like graphs and fingerprints as well as actions like the flow of functional or experimental information, deep generative models can theoretically produce unique chemical and biological structures with the desired features. Deep generative models that are imaginative can greatly advance algorithm development and use in drug discovery. Deep generative models would present cutting-edge technologies that could revolutionize an informatics understanding of biology, illness, and therapies in this "big data" era. It takes a lot of work to create a novel medicine that meets all the requirements for on-target potency, selectivity in relation to off-targets, physical qualities, and other chemistry and biology parameters. Chemists must choose and experimentally validate candidate compounds from a broad chemical space using the conventional procedures, which are unsuccessful. Deep generative models have gained popularity because they can quickly and cheaply produce new bioactive and synthesizable compounds. Several widely used chemical and bioinformatics

*Editorial Office, Journal of Pharmacology and Medicinal Chemistry, Windsor, Berkshire, England*

*Correspondence: Christopher Hopkins, Editorial Office, Journal of Pharmacology and Medicinal Chemistry, Windsor, Berkshire, England, e-mail jpharmacology@theresearchpub.com*

databases that give the drug development community access to labeled and unlabeled data for training, validating, and testing deep generative models. The internal proprietary libraries of pharmaceutical corporations range from 2 to 3 million chemicals and include data from prior attempts at medication discovery. For in silico screening, the ZINC database gathered roughly 2 billion buyable, readily available "drug-like" chemicals from the public domain. Due to its enormous scale, it is also effective for pre-training generative models by teaching molecular patterns. It is particularly interesting to study bioactive molecules, such as those in the manually curated ChEMBL database, which contains about 1.5 million actual bioactive chemicals, each of which has at least one experimental bioactivity measurement. They may be employed to train models to produce molecules with particular characteristics. The majority of organic compounds (166.4 billion) with up to 17 heavy atoms of C, N, O, S, and halogens are listed in the GDB-17 database. This covers a lot of the smaller, common lead compounds as well as medications with lower molecular weight. Chemoinformatics techniques and expert-system-style rules have been used to identify billions of synthesizable molecules in extremely large chemical databases like Enamine and REALdb. The option to train models with broader applicability is provided by these extremely big databases. A number of macromolecular databases, including the PDB, provide richer data for generative model training in macromolecule design in addition to small-molecule resources. In order to process human language, generative neural networks must include Recurrent Neural Networks (RNNs). They have proved effective in automating NLP computer code production and musical composition, and they are valuable for modeling systems that have a sequential or temporal component. Human language and the language of molecules, like SMILES, are similar. Thus, using RNNs to generate molecules based on sequential representation is natural. The two networks that make up an Autoencoder (AE) are the encoder, which is taught to map input data into a low-dimensional latent vector, and the decoder, which is trained to map the latent vector back into the input data. The original AE copies the input to produce a latent space. Variational AE (VAE) regularizes the latent space by substituting latent space distributions for latent space points in order to prevent overfitting and discontinuities in the original AE. In a ground-breaking study, VAE was used to generate molecules, launching a brand-new approach to de novo drug discovery. A molecule is represented as an explicit probability distribution over latent space because the latent vectors are forced to follow a probability distribution (often Gaussian distribution). When the encoder and decoder are trained together, the output must recreate the probability distribution of the training samples. The basic objective of learning disentangled representations for VAE, which aim to make each latent variable of the latent vector encode an independent attribute or aspect of data, has recently attracted significant attention. If disentangled VAE is effectively introduced for molecular generation, a molecular property can be changed by altering the latent variables related to that property without altering other qualities.