

Revolutionizing Pharmaceuticals: Applications and Potential of Generative Artificial Intelligence in Drug Discovery

Ángel Canal-Alonso¹, Noelia Egido¹, Pedro Jiménez¹, Javier Prieto¹, Juan Manuel Corchado¹

¹Department of Bioinformatics and Computational Biology, AIR Institute, Carbajosa de la Sagrada, Spain

Email: acanal@air-institute.com

Abstract

Artificial intelligence (AI) has emerged as a transformative tool in the pharmaceutical industry, revolutionizing the traditional drug discovery and development process. Through advanced generative techniques, such as Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs), the exploration and design of novel and viable therapeutic molecules has been enhanced. Additionally, AI facilitates the optimization of these molecules by guaranteeing desirable properties and accelerates the identification of therapeutic targets through deep analysis of biomedical and genomic data sets. One of the most significant advances has been drug repurposing, where AI unlocks the hidden potential of known drugs for new therapeutic indications.

Keywords: Generative Artificial Intelligence, Drug Design, Computational Models

1. Introduction

Generative artificial intelligence (AI) is a prominent subspecialty within the broad field of artificial intelligence. This focuses primarily on the design and development of computational models capable of producing data that emulate those present in a specific training set. In other words, rather than simply analyzing and processing data, generative AI strives to “create” or “generate” data that closely reflects the essence of the information it has been trained on.

The heart of these models lies in advanced machine learning techniques. A prominent example of this is generative neural networks, which, through mathematical structures and sophisticated algorithms, can fabricate new examples of data that follow the same distribution as the original data. These systems, by learning underlying patterns and structures in

large amounts of data, can produce results that, in many cases, are virtually indistinguishable from authentic data.

In sectors such as the pharmaceutical industry, the potential of generative AI has been recognized and is being harnessed in revolutionary ways. With the constant challenge of discovering and developing more effective and safe medicines in shorter times, generative AI presents itself as an invaluable tool. It can simulate molecular interaction, predict the efficacy of new substances and dramatically accelerate the research phase in the creation of new drugs. Furthermore, beyond the direct development of medicines, generative AI contributes to optimizing processes, personalizing treatments and improving efficiency at various points in the pharmaceutical value chain. For example, in the production and distribution of medicines, it can help anticipate demands or foresee ideal production scenarios,

Generative neural networks, often referred to as Generative Adversarial Networks (GANs), represent one of the most exciting advances in the field of machine learning in the last decade. These networks, as their name suggests, focus on data generation, and their unique design allows them to produce information that is strikingly similar to real data distributions.

These GANs operate through a binary structure: the generator and the discriminator. The interaction between these two components is essentially a game of cat and mouse, and it is precisely this interaction that allows GANs to achieve such impressive results.

The generator, as its name suggests, is responsible for generating data. It starts by taking an initial input that is usually random, often based on a simple noise distribution. From this input, the generator uses its neural architecture to transform that noise into data that attempts to emulate or imitate real examples. Initially, this generated data may not look much like reality, but with time and proper training, the accuracy improves significantly.

On the other hand, the discriminator functions as a judge. His task is to evaluate the presented data and decide whether it comes from the real training set or whether it has been generated by the generator. At the beginning, when the generator is still in the early stages of its training, this task is relatively simple for the discriminator. However, as the generator becomes more skilled, the task becomes more complicated.

The true power of GANs comes from the feedback between these two components. When the discriminator correctly identifies a generated piece of data as false, it sends feedback to the generator, telling it where and how it fell short. In turn, the generator uses this information to adjust and improve. This iteration continues in a loop, in which the generator constantly tries to outperform the discriminator, and the discriminator strives not to be fooled.

Over time, this competition leads the generator to produce data of such quality that the discriminator has difficulty distinguishing between what is real and what is generated. This iterative and competitive process, although simple in theory, has led to astonishing advances in fields as varied as image generation, music creation, data simulation, and much more. GANs, with their unique and powerful structure, have revolutionized the ability of machines to create and understand complex data.

Generative artificial intelligence (AI) has brought about a revolution in the way we approach and use large data sets, especially when it comes to unsupervised learning. Unlike

supervised learning, where models are trained using data with clear labels, unsupervised learning works with data that does not have such labels. Instead of being directed toward a specific answer, the model is trained to discover the underlying structures and patterns in the data on its own.

This self-directed form of learning is especially relevant and valuable in fields where you have access to massive amounts of unlabeled data. And the pharmaceutical industry is a clear example of this. Scientists and researchers in this field frequently work with vast data sets related to molecules, chemical structures, genetic sequences, biological activity profiles, and more. Many times, this data does not come with clear labels or specific definitions, making traditional or manual processing a Herculean, if not impossible, task.

This is where generative AI using unsupervised learning comes into play. These systems, by not depending on predefined labels, have the ability to immerse themselves in these oceans of data and, through advanced algorithms, detect patterns, relationships and connections that could go unnoticed by humans. In doing so, these AIs can identify, for example, how certain molecular structures relate to specific biological activities or how different combinations of molecules could result in effective therapeutic compounds.

Discovering hidden patterns and identifying potentially significant relationships through unsupervised learning has profound implications for the pharmaceutical industry. It can accelerate the drug discovery process by reducing the search space, aligning candidate compounds with specific diseases more quickly, and providing a deeper understanding of molecular interactions. Furthermore, by making the research process more efficient, associated costs can also be reduced and the arrival of innovative therapeutic solutions to the market can be accelerated.

2. Common architectures

There are several common generative techniques used in the pharmaceutical industry, including:

Generative Adversarial Networks:

Generative Adversarial Networks, known by their acronym in English as GANs, represent one of the most notable advances in the field of artificial intelligence in recent years. They were introduced to the world by Ian Goodfellow and his team in 2014, and since then, they have revolutionized fields as diverse as image creation, sound design, data simulation, and much more.

What is truly innovative about GANs is the way they are structured and how they operate. As indicated, these networks are composed of two distinct but interconnected models: the generator and the discriminator. Together, they participate in a kind of competitive "game", where each model seeks to outdo the other in a constant tug-of-war.

The generator has the task of creating or generating data. Initially, you can start with a simple random input, often derived from a noise distribution. Through its neural network structure, this model strives to transform that input into a sample that resembles real data. However, especially in the early stages, what it produces is not perfect and can be easily discernible from the actual data.

This is where the discriminator comes into play. Acting as a critic or judge, your role is to evaluate the samples and determine if they are real (coming from the original data set) or if they are fake (generated by the generator). At the beginning of training, when the generator is still learning, it is easier for the discriminator to make this distinction. But as the generator improves, the discriminator's job becomes more challenging.

What makes GANs so efficient is the feedback dynamics between these two models. Every time the discriminator correctly identifies a generated sample, it provides information to the generator on how to improve. This cycle of generation, discrimination, and feedback continues until the generator reaches a level where it can produce samples so convincing that the discriminator has difficulty distinguishing them from real data.

Since their introduction in 2014, GANs have undergone a series of improvements and variations, giving rise to a wide range of applications and derived models. They have been used in art creation, music generation, 3D model creation, image super-resolution, and many other fields. The ability of GANs to generate realistic data from random inputs has proven to be a powerful and versatile tool in the world of machine learning and artificial intelligence.

Variational Autoencoders:

Variational Autoencoders, or VAEs, represent a fascinating evolution in the field of generative models, and although they share similarities with GANs in their ability to generate data, their underlying approach and their operating mechanics are different.

Let's start by unpacking a little what an "autoencoder" means. In machine learning, an autoencoder is a neural

network used to learn encoded (or compressed) representations of input data, so that these data can be reconstructed from these representations with the least possible loss of information. Basically, it takes data, compresses it into a latent representation, and then attempts to decompress or "decode" that representation to get a reconstructed version of the original data.

Now, what makes an autoencoder "variational"? The variation is introduced in how the latent representation is modeled. Instead of learning a fixed encoded representation for each input data, a VAE learns the parameters of a probability distribution for the latent representation. This introduces some randomness into the process, meaning that every time we encode and then decode an input, the result may vary slightly.

The magic behind VAEs lies in their ability to generate new data. Once trained, we can sample random values from the latent distribution and pass them through the decoder part of the VAE to obtain new data. These generated data will be consistent with the overall distribution of the training data.

Now, in relation to the pharmaceutical industry, VAEs are presented as potentially revolutionary tools. Imagine having vast databases of molecules or compounds and wanting to explore slightly different variants or completely new compounds that have not yet been synthesized. By training a VAE on such data, researchers can sample from the latent space to generate novel molecular structures that could have desirable pharmacological properties.

What distinguishes VAEs from GANs, in addition to their internal mechanism, is their theoretical basis. While GANs are based on a competition between two networks, VAEs use the concept of variational inference, which is a method for approximating complex probability distributions through simpler distributions. This makes them particularly suitable for working in scenarios where you want to have more explicit control over the distributions of the data generated or when you have an interest in the latent space itself.

Generative Flow Models:

Generative flow models, known simply as "flows", represent an advanced and promising branch within the field of generative models. Unlike the previously mentioned techniques, such as GANs and VAEs, flows are characterized by operating directly on transformations of probability distributions, allowing more explicit and direct control in the data generation process.

The essence of a generative flow model lies in its ability to transform a simple probability distribution, such as a standard Gaussian distribution, into a more complex and specific probability distribution that resembles the distribution of the data with which it has been analyzed. This is achieved through a series of invertible and differentiable transformations that, together, are called "flows." These transformations retain the ability to sample from the simple distribution and at the same time adjust it to match the desired distribution.

A distinctive advantage of flows is their ability to accurately and efficiently calculate the probability density of the generated samples. In practical terms, this means that you can obtain not only a generated sample, but also a measure of how "probable" or "coherent" that sample is with respect to the training data.

Within the context of the pharmaceutical industry, this feature is especially valuable. When generating new compounds or molecules, it is not only essential to have a coherent chemical structure, but also to understand the probability that such a structure can arise given existing data distributions. Thus, generative flow models can offer a double benefit: on the one hand, they generate high-quality data and, on the other, they provide a metric of the quality of that data.

Furthermore, the deterministic structure and the ability to directly manipulate probability distributions make flow models more interpretable than some other generative techniques. This can be crucial in pharmaceutical research, where the interpretation and understanding of the data generated can be as important as the generation itself.

Transformers.

The Transformers architecture has triggered a true revolution in the field of deep learning and artificial intelligence. Its ability to capture contextual relationships between words or tokens, regardless of their relative position in a sequence, has catapulted its relevance, especially in natural language processing (NLP).

Originally conceived to address challenges in NLP, as highlighted in the influential work of Vaswani et al. In 2017, Transformers have expanded beyond these initial applications. Its distinctive mechanism, known as "self-regressive attention," allows the network to assign different relevance weights to each word or token based on its context, which is critical for understanding meaning and semantics in human languages. It is this focus on context and attention that has allowed Transformers to outperform previous neural

architectures, such as convolutional (CNNs) and recurrent neural architectures (RNNs), in numerous NLP tasks.

In addition to traditional NLP tasks such as machine translation and text generation, the flexibility of Transformers architecture has made them suitable for a variety of non-language applications. In the context of the pharmaceutical industry, where data structures can be as complex and hierarchical as human language, Transformers are finding a growing niche. Molecules and proteins, for example, have intricate structures and relationships between their components that can be compared metaphorically to the relationships between words in a text.

Applying Transformers in the pharmaceutical field involves treating molecular structures and protein sequences as "texts", where atoms, bonds and amino acids can be seen as "words" or "tokens". By doing so, it is possible to train models that understand and generate molecules with desirable properties or that predict how a given molecule will interact with a specific biological target.

Transformers' ability to generate high-quality, consistent data has proven invaluable for molecule design in pharmaceutical research. These models can, for example, be trained with databases of existing molecules and then generate novel drug candidates that could have beneficial pharmacological properties.

In conclusion, although Transformers emerged from the world of natural language processing, their versatility and power have been applied in numerous fields, including pharmaceutical design and development. Their influence is laying the foundation for a new era in generative artificial intelligence and its interaction with cutting-edge science.

3. Use of Generative AI for drug discovery

The traditional drug discovery process is long, expensive, and involves numerous trials and errors. In this context, generative artificial intelligence has proven its value by streamlining and optimizing several crucial stages of new drug discovery. Below are some of the most notable applications of generative AI in this field.

Molecule design

The traditional drug design and discovery process is extremely complex, laborious and expensive. Historically, it required extensive trial and error to identify and optimize compounds with the potential to become effective medicines. However, with the rise of artificial intelligence (AI) and specifically generative AI, this dynamic is undergoing a radical change.

Molecule design is a discipline that merges chemistry with bioinformatics, seeking to design chemical structures with specific pharmacological properties. Given the immensity of chemical space—that is, the totality of possible molecules—manual or even semi-automated exploration of this space is a titanic task. This is where generative AI comes into play as a revolutionary tool.

Using advanced machine learning techniques, such as GANs and VAEs, it is possible to model and sample this vast chemical space. GANs, with their generator and discriminator structure, can produce molecules that are indistinguishable from real molecules, while VAEs allow a structured exploration of chemical space, sampling and decoding points in a latent space to generate new molecular structures.

What makes generative AI especially powerful in this context is its ability to generate molecules that are not only novel, but also chemically and pharmacologically viable. This means that the proposed molecules are not only structurally feasible, but also present desired properties, such as high affinity for a specific biological receptor, low toxicity, and feasibility of synthesis in the laboratory.

Furthermore, by integrating existing databases on molecular properties, biological activity and toxicity, generative models can be trained to take these factors into account when proposing new molecules. This can lead to a reduction in the number of compounds that need to be synthesized and tested experimentally, saving time and resources.

In practical terms, the adoption of generative AI in molecule design is leading to a new era in drug discovery. These advances allow researchers and pharmaceutical companies to explore previously inaccessible regions of chemical space, identify drug candidates more quickly, and ultimately bring innovative treatments to market in a more efficient and cost-effective manner.

Optimization of molecular properties

The drug discovery process is not only about identifying molecules that may have a desired biological activity, but also about ensuring that these molecules are suitable for administration and therapeutic use in humans. A molecule can have a potent effect on a biological target, but if it is unstable, insoluble in water, or shows toxicity, it will not be viable as a drug. Therefore, the optimization of molecular properties becomes essential.

Optimization is traditionally performed through an iterative process that involves modifying the chemical structure of a molecule and then evaluating the resulting properties. This process can be long and expensive. However, with the inclusion of generative AI, we can navigate this process much more strategically and efficiently.

Generative AI, particularly when combined with powerful architectures like Transformers, can quickly explore variations of a molecule, predict how these modifications will affect its properties, and suggest the most promising versions of the molecule. For example, if a candidate molecule shows activity but has solubility issues, AI can suggest structural modifications that could improve this solubility without compromising activity.

Transformers, with their ability to understand and model complex relationships, are ideal for this type of task. They can capture the subtleties and intricate relationships in molecular structures and translate that understanding into molecular design suggestions. This ability to "understand" and "reason" about chemical structures at such a detailed level surpasses many previous techniques.

Furthermore, these generative AI tools not only optimize based on a single criterion, but can take into account multiple objectives simultaneously. This is crucial in drug design, where properties such as activity, selectivity, toxicity and solubility are often interrelated and can compromise with each other.

In summary, the adoption of generative AI, and in particular architectures such as Transformers, in the molecular optimization process is facilitating an era of more precise, faster and cost-efficient drug design. With these tools at their disposal, researchers can directly target molecules that are not only active but also suitable for clinical development and eventual administration to patients.

Identification of therapeutic objectives

The drug discovery process begins with the identification of a suitable therapeutic target, that is, a molecule or pathway in the body that, if modified in some way, could lead to a therapeutic response. While identifying these targets is crucial, it is not a simple task due to the inherent complexity of biological systems and the vast amount of biomedical and genomic data available.

Historically, the identification of therapeutic targets was based on laboratory experiments and a general understanding of the biology of a disease. But, in the era of systems biology

and genomics, we are inundated with data that may be unexplored or underutilized due to human limitations in analyzing large volumes of information.

This is where generative AI offers revolutionary potential. By applying advanced machine learning techniques, AI can analyze and synthesize large data sets, identifying patterns and connections that might go unnoticed by a human researcher.

Natural language processing (NLP), for example, is a powerful tool in this context. It can be used to analyze scientific literature, patents and other related texts to identify mentions of genes, proteins or pathways that are associated with specific diseases. By correlating this information with genomic and biomedical databases, AI can highlight targets potentially relevant to a particular disease.

Similarly, biological network analysis allows AI to explore how different molecules and pathways interact with each other in the context of a cell or tissue. This is essential to identify points of intervention that could have maximum therapeutic impact with minimum side effects.

Another advantage of generative AI in this context is its ability to integrate and analyze heterogeneous data, from transcriptomic and proteomic data to clinical and epidemiological information. By combining all of these sources, AI can generate a more complete and accurate picture of a potential therapeutic target.

Ultimately, by automating and optimizing the identification of therapeutic targets, generative AI not only accelerates this crucial initial process in drug discovery, but also increases the likelihood of success in later stages of drug development, leading to therapeutics, more effective and safer products to the market in a shorter period of time.

“Repurposing” of drugs

The process of discovering and developing new drugs is expensive, risky and can take many years. Given this intensive investment in time and resources, drug repurposing—which involves finding new therapeutic uses for existing medications—is presented as an attractive and efficient alternative. Because these drugs have been through clinical testing and their safety profiles are known, repurposing can significantly reduce the costs and time associated with traditional drug development.

Generative AI plays a fundamental role in the modernization and efficiency of this process. Instead of relying exclusively on trial and error or serendipity, AI can

systematically scan and analyze vast amounts of information to identify promising candidates for repurposing.

For example, by extracting information from large databases containing information on drug-protein interactions, generative AI can discover previously unrecognized modes of action for existing drugs. Additionally, by analyzing scientific literature and patient data using natural language processing techniques, AI can identify correlations and patterns between diseases and treatments that are not evident to the naked eye.

Another approach is the analysis of molecular signatures. AI can compare the molecular responses induced by different drugs with the signatures of various diseases. If a signature induced by a drug is opposite to the signature of a disease, that drug could be a candidate to treat that disease.

Importantly, while generative AI provides valuable clues and predictions, these must be validated through experiments and clinical trials to confirm their efficacy and safety in proposed new indications.

Using generative AI in drug repurposing not only accelerates the drug development process but also expands the therapeutic potential of existing drugs. In a world where unmet medical needs continue to be a challenge, this AI-backed strategy holds a hopeful promise for patients and healthcare professionals.

4. Future work

The adoption and success of generative artificial intelligence in the pharmaceutical industry to date marks just the beginning of what could be a series of significant advances in medicine and therapeutics. As we move forward, it is essential to explore and delve in various directions to maximize the potential of AI in this field.

First, while current AI tools have proven effective in exploring chemical space, it is essential to develop more advanced algorithms that can address the complexity inherent in human biology and molecular-molecular interactions. This will allow for more accurate prediction of the pharmacokinetic and pharmacodynamic properties of drug candidates.

Furthermore, the integration of multiple data sources, such as genomics, proteomics, metabolomics, and transcriptomics, can provide a holistic and detailed view of the impact of a drug candidate on a biological system. Working on algorithms that can efficiently handle and analyze these large, heterogeneous data sets will be essential.

The interpretability and transparency of AI models is also an area that requires attention. As AI makes more complex decisions in drug design and optimization, it is critical that scientists and regulators understand how these decisions are made to ensure the safety and effectiveness of proposed treatments.

Finally, as AI becomes a standard tool in pharmaceutical research, it will be crucial to develop regulatory and ethical frameworks that ensure its appropriate and safe use. Collaboration between researchers, industry and regulatory bodies will be essential to establish guidelines that support innovation without compromising patient safety.

In conclusion, although we have witnessed notable advances thanks to generative AI in pharmaceuticals, the path towards its full and optimal integration is still being charted. Investment in research and development, as well as interdisciplinary collaboration, will be key to ensuring that this technology reaches its full potential to benefit human health.

5. Conclusions

The revolution that artificial intelligence (AI) has introduced in multiple areas has had a transformative impact on the pharmaceutical sector. Its application in drug discovery and development has promised and has already begun to deliver a more agile and precise era of pharmaceutical innovation. Several aspects of this transformation deserve special attention.

First, by addressing the fundamental challenge of molecular design, generative AI has proven to be a powerful tool. By using architectures such as Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs), the ability to explore vast chemical spaces in search of molecules that are not only novel but also viable from a therapeutic point of view has been enhanced. This exploration, previously tedious and often based on trial and error, has been significantly accelerated.

The optimization of these molecules to guarantee desirable properties, such as solubility, stability and selectivity, has been another field of action of generative AI. Using advanced techniques, it is possible to refine existing molecules or propose new structures that meet precise criteria, facilitating the transition from a promising compound to a viable clinical candidate.

Furthermore, at the very beginning of the drug discovery process, the identification of suitable therapeutic targets is essential. In this area, AI has proven to be exceptionally

useful, providing deeper and more systematic analysis of large biomedical and genomic data sets. This ability to discern complex patterns and hidden connections has led to more informed and evidence-based discoveries in the realm of therapeutic goals.

However, perhaps one of the most promising strategies where AI has left its mark is drug repurposing. The reuse of existing drugs for new therapeutic indications is an efficient and profitable way to address diseases without adequate therapeutic solutions. Here, AI not only proposes candidates for repurposing based on existing data but also uncovers previously unknown relationships and modes of action, unlocking the hidden therapeutic potential of already known drugs.

Overall, while traditional pharmaceutical research has been, and continues to be, a cornerstone in the advancement of medicine, the adoption and adaptation of tools based on generative AI represent a paradigmatic shift. These tools are poised to accelerate the pace of discovery, reduce costs, increase precision and, most importantly, improve the chances of success in the search for more effective and safer treatments. As we continue to navigate this era of digital transformation, it is evident that deep integration of AI into the pharmaceutical industry is not only desirable, but essential to meeting the medical challenges of the 21st century.

References

- Garcia-Retuerta D, Canal-Alonso A, Casado-Vara R, Rey AM, Panuccio G, Corchado JM. Bidirectional-Pass Algorithm for Interictal Event Detection. In Practical Applications of Computational Biology & Bioinformatics, 14th International Conference (PACBB 2020). PACBB 2020. Advances in Intelligent Systems and Computing, vol 1240. Springer, Cham. https://doi.org/10.1007/978-3-030-54568-0_20
- Castillo Ossa LF, Chamoso P, Arango-López J, Pinto-Santos F, Isaza GA, Santa-Cruz-González C, Ceballos-Marquez A, Hernández G, Corchado JM. A Hybrid Model for COVID-19 Monitoring and Prediction. *Electronics*. 2021; 10(7):799. <https://doi.org/10.3390/electronics10070799>
- Intelligent Platform Based on Smart PPE for Safety in Workplaces. Márquez-Sánchez S, Campero-Jurado I, Herrera-Santos J, Rodríguez S, Corchado JM. *Sensors (Basel)*. 2021 Jul 7;21(14):4652 <https://doi.org/10.3390/s21144652>
- A. Canal-Alonso, R. Casado-Vara and J. Manuel Corchado, "An affordable implantable VNS for use in animal research,"

2020 27th IEEE International Conference on Electronics, Circuits and Systems (ICECS), 2020, pp. 1-4, doi: 10.1109/ICECS49266.2020.9294958

An Agent-Based Clustering Approach for Gene Selection in Gene Expression Microarray. Ramos J, Castellanos-Garzón JA, González-Briones A, de Paz JF, Corchado JM. *Interdiscip Sci.* 2017 Mar;9(1):1-13

DOI 10.1007/s12539-017-0219-6

Acknowledgments

This study has been funded by the AIR Genomics project (with file number CCTT3/20/SA/0003), through the call 2020 R&D PROJECTS ORIENTED TO THE EXCELLENCE AND COMPETITIVE IMPROVEMENT OF THE CCTT by the Institute of Business Competitiveness of Castilla y León and FEDER funds.