

# Ask the Experts: The impact of artificial intelligence in drug discovery

Artificial intelligence (AI) has become more common, both in our research labs and in our homes, but what are the limitations of AI?

We turn to Anne Carpenter (Broad Institute; MA, USA), Wengong Jin (Eric and Wendy Schmidt Center; MA, USA), and Jürgen Bajorath (University of Bonn; Germany) to answer our questions about developing computational techniques for drug discovery, the challenges of doing so, and how this technology might evolve in the future.

## Contents

---

• What is the difference between AI, machine learning and deep learning? .....	2
• In what ways can AI be used to accelerate drug discovery?.....	3
• What have been some significant advancements or successes of AI in drug discovery?.....	4
• How is AI being used within your own research?.....	5
• What do you think is a common misconception about using AI in drug discovery?...	6
• What are the challenges of developing AI for drug discovery?.....	6
• What are the current limitations of using AI in drug discovery?.....	7
• How do you think AI will evolve in the next decade to accelerate drug discovery?...	8
• Meet the Experts.....	10

# What is the difference between AI, machine learning and deep learning?

Anne Carpenter

These terms can be confusing because some envelop the others, and some have both a technical meaning and an everyday meaning. Simply speaking, in machine learning (ML) you aim to teach a computer to answer questions correctly by providing it with examples (either examples with the correct answers, in supervised learning, or just examples of the data in unsupervised learning). The computer aims to discover general rules rather than just memorizing answers. ML can be trained to answer questions like, "Where are the nuclei in this image?" or "Where are the transcription factor binding sites in this genome sequence?" or "What groups of similar samples exist in this dataset?"

To understand deep learning (DL), it helps to know that in most ML applications to date, specific features were extracted intentionally from the data in the hopes that those features would make it easier for the computer to learn correct answers. For example, we design features in images relating to the texture, shape, and size of cytoplasmic staining to try to detect whether a cell is metastatic or not.

However, DL is a type of ML where instead, you feed the raw data to the computer, usually in huge quantities, and let it sort out how to best extract features from the data in order to make its decisions. The decision-making system has many internal layers, which sparked the name 'deep' learning. In my metastatic cell example, we would just give the system the raw image pixel data and let the system figure out how to distinguish metastatic cells by giving it many examples.

Now, AI: some use the term artificial intelligence to refer to any computer system that can make good decisions, whereas others use it to refer just to a branch of ML where the computer is forced to understand how it can learn generally, rather than being trained just for a specific task.

Jürgen Bajorath

ML is a sub-discipline of AI and DL uses deep neural network (DNN) architectures and is a sub-discipline of ML.

# In what ways can AI be used to accelerate drug discovery?

## Anne Carpenter

This is such an exciting time to be working at the interface of computer science and drug discovery because there are so many applications! Many individual steps of drug discovery can be accelerated using AI.

For example, you can train systems to predict a given compound's activity in an assay based on its chemical structure or other pre-existing information about the compound. You can train a system to sift through millions of images of cells treated with compounds to identify a favorable phenotype. You can predict the structure of a protein involved in a disease so that compounds can be designed to fit into them better. You can even test millions of chemical structures virtually to assess how they bind with the protein. Systems for these and many more tasks are not perfect, but they can assist experts in making better and faster decisions.

## Wengong Jin

Molecular screening is a crucial step in drug discovery, where a chemist puts a library of existing compounds into a biological assay to measure their biological properties, such as potency, toxicity and solubility. The number of chemicals that could be potential drug candidates is estimated to be at least  $10^{60}$  – these are all the molecules that obey Lipinski's rule-of-five for oral bioavailability – creating a major bottleneck in screening for drug candidates.

Standard high-throughput screening facilities in the pharmaceutical industry can only test around  $10^5$  compounds per day. It is, therefore, crucial to restrict the size of compound libraries to make the screening time and associated costs feasible. We seek to accelerate and automate drug discovery using AI.

Previous screening efforts in the pharmaceutical industry have generated many datasets of molecules with labeled properties. This allows us to build molecular property models that can predict the properties of a compound without testing it in a wet lab. We can then use these models to virtually screen a much larger collection of molecules at a much faster speed ( $10^8$  compounds/day) than is possible with current high-throughput screening facilities in a wet lab.

## Jürgen Bajorath

The hope is that AI approaches will further expand the currently charted chemical/target space and accelerate discovery paths from targets and novel chemical entities to drug candidates.

AI enterprises engaged in drug discovery already claim such accomplishments on a case-by-case basis. However, a word of caution is advisable since there is typically a gap between claims, promotional efforts (for example, for fundraising), and the scientific reality when it comes to pushing 'new' technologies in drug discovery. In this context, it should be noted that ML has a long history in pharmaceutical research and that DL represents an extension of this framework, rather than a truly novel approach.

## What have been some significant advancements or successes of AI in drug discovery?

### Anne Carpenter

I serve on the scientific advisory board of a company called Recursion, which uses ML to identify changes in cell morphology that occur when genes associated with disorders are perturbed. The team then screen compounds to identify those that could reverse disease-associated changes. Using computers to analyze images makes these decisions fast and objective. They now have four candidate therapeutics entering clinical trials!

### Wengong Jin

In 2020, whilst I was at MIT (MA, USA), we successfully used AI to discover a new antibiotic called Halicin. We did this by training a DNN to become capable of predicting molecules with antibacterial activity. We performed predictions on multiple chemical libraries and discovered Halicin, a compound that is structurally divergent from conventional antibiotics and displays bactericidal activity against a wide spectrum of pathogens including *Mycobacterium tuberculosis* and carbapenem-resistant *Enterobacteriaceae*.

Halicin also effectively treated *Clostridioides difficile* and pan-resistant *Acinetobacter baumannii* infections in mice. This work highlights the utility of DL approaches to expand our antibiotic arsenal through the discovery of structurally distinct antibacterial molecules. This discovery was published in *Cell* and received significant attention because there is an urgent need to discover new antibiotics due to the rapid emergence of antibiotic-resistant bacteria.

## Jürgen Bajorath

Currently, AI in drug discovery mostly refers to DL and robotics, while the adaptation of other AI sub-disciplines is still at very early stages. One of the areas where DL has recently made a substantial impact is in computer-aided synthesis planning and prediction, at least at the methodological level. However, many medicinal chemists attest to the fact that these DL-driven advances are yet to be made practically applicable in their day-to-day efforts, aside from raising awareness of what this technology can do. Time and substantial efforts will be required until AI/DL tools measurably impact the practice of drug discovery on a larger scale.

## How is AI being used within your own research?

### Anne Carpenter

We are teaching the computer to see things that humans cannot see in images. For example, by eye, humans cannot distinguish cells with a certain type of leukemia from those without, so biomarkers were developed that could be detected by fluorescence flow cytometry. We recently used DL to teach the computer to identify those leukemic cells based on just unstained microscopy images, without any biomarker labels, and it succeeded!

### Wengong Jin

I am currently using AI to search for synergistic drug combinations to treat COVID-19. Drug combinations make promising therapeutic candidates for COVID-19, but the lack of high-quality training data makes it difficult for DL to predict drug synergy accurately.

To address this challenge, I proposed a novel DL model called ComboNet, which jointly models drug-target interaction and drug synergy. Together with the National Center for Advancing Translational Sciences (MA, USA), we discovered two novel drug combinations (remdesivir and reserpine; remdesivir and IQ-1S) with strong synergy. This work was published in *PNAS* in 2021 and we are currently applying this model to find effective drug combinations for pancreatic cancer.

### Jürgen Bajorath

Our research largely focuses on computer-aided medicinal chemistry and chemoinformatics. Like other groups in this area, we have been using ML for molecular property predictions and other applications for many years. Furthermore, we have also developed ML approaches for a number of specific tasks such as predicting activity cliffs (structurally similar compounds that are active against the same target but with large differences in potency) or compound-target screening matrices.



In recent years, I have become increasingly interested in better understanding ML predictions, their successes, and failures (or, in more colorful terms, shedding light on the 'black box' of ML, when even the designers of a computer model cannot explain how a certain decision is made). This is also referred to as 'explainable AI' (XAI).

XAI refers to methods that allows humans to comprehend the results outputted by ML algorithms. Notably, one of the attractions of DL is that DNNs enable us to tackle problems that are difficult, if not impossible, to address using standard ML approaches such as molecular image-based predictions or chemical representation learning. This is another major driver for increasingly investigating DL in our research environment.

## What do you think is a common misconception about using AI in drug discovery?

Jürgen Bajorath

Firstly, it is often not sufficiently understood what AI is – and what it is not. We are still far away from a situation where computers make autonomous decisions beyond human reasoning, at least in pharmaceutical research. DL is data-driven, statistical in nature, and far from being some form of 'magic' for unsolved problems in drug discovery, such as high attrition rates.

Secondly, high expectations that AI might 'revolutionize' the drug discovery process are on rather fragile grounds. No single scientific approach or technology has ever come close to revolutionizing drug discovery and there are good reasons to anticipate that this will also apply to AI. Hence, in light of the drug discovery history, arriving at a better general understanding of current AI approaches, their opportunities and limitations, would be beneficial for pharma environments and help to avoid unrealistic expectations.

## What are the challenges of developing AI for drug discovery?

Anne Carpenter

It's fairly easy to achieve successful results for a supervised ML problem if you try enough parameters or architectures and test it on only a small sample that is very similar to what you've it trained on. The challenge is to create something that works reliably in the real world, and that takes a serious investment in creating the training and testing data to be sure that you are not fooling yourself with a system that has just memorized the correct answers for a small dataset.

## Wengong Jin

The major challenge of developing AI for drug discovery is data scarcity and bias, as training data is usually limited in molecular property prediction, or is otherwise biased. Additionally, molecular assays used for learning property predictors involve many sources of spurious correlations, as a result of the choice of chemical libraries, batch effects, or measurement biases, for example. Therefore, effective molecular property prediction requires that models generalize beyond the chemical space of training examples and avoid learning spurious correlations introduced by these biases.

It is also challenging to design proper evaluation protocols to measure the generalization power of a method when applied to a new chemical space, as is common in drug discovery.

## Jürgen Bajorath

Unlike other fields where AI/DL has made a strong impact, drug discovery is overall not a data-rich discipline. The use of limited amounts of mostly structured data does not play into the strengths of 'data-hungry' DL approaches. Consequently, consistent improvements of DL predictions over other ML approaches are not expected across typical applications such as compound activity or property predictions and are currently not observed.

In drug discovery settings, it will be important to identify applications where DL is most likely to outperform standard ML approaches (for example, image-based analysis of high-content assays) and concentrate on novel applications that are essentially enabled through DL (such as advanced synthesis design). In addition to data constraints, it should also be taken into consideration that drug discovery is a highly interdisciplinary process with intrinsic scientific heterogeneity, making it rather unlikely that 'one-size-fits-all' AI systems will be easy to conceptualize and implement.

## What are some of the current limitations of using AI in drug discovery?

### Anne Carpenter

One of the biggest challenges I see is in predicting the toxicity of compounds. Solving this problem would have a HUGE impact on the pharmaceutical industry but it's very challenging to design AI solutions for it. For example, if I invented a new AI-based tool that could tell you whether a given chemical structure would be toxic to humans, how could I prove it works? We can't give lots of different compounds to humans outside of clinical trials, and there are very few trials of new compounds each year to validate my system.



I could train my system to predict the outcomes of toxicity testing on animals, but we know that animal results are not entirely consistent with human results (although are better than nothing). We could test the system against past clinical trials, but most likely that is the data I used to train my system, so it might have just memorized the right answers. So, the very small dataset of human toxicity data is a major challenge.

### Jürgen Bajorath

In addition to general limitations resulting from data sparseness, the black box character of AI/DL is another important issue. Drug discovery practitioners are typically reluctant to rely on predictions that cannot be understood in chemical or biological terms, which works against the acceptance of black box approaches for practical applications. This emphasizes the need for XAI methods to rationalize predictions and communicate them in an intuitive manner.

Since operating in discovery project teams typically requires multi-tasking and working under time pressure, ease-of-use and robustness of new computational methods and tools are essential for using them in practical applications and for making progress. While developing consistently accurate predictive models is a formidable challenge, transforming expert domain models into widely accessible tools presents another challenge of similar magnitude.

## How do you think AI will evolve in the next decade to accelerate drug discovery?

### Anne Carpenter

ML will be incorporated more seriously at each step in the pipeline, providing assistance to experts and making their work more efficient. On top of this, I imagine we will see improvements in generative ML systems. So, instead of telling you whether a proposed compound is likely to be effective, this can instead generate a structure from scratch that is predicted to have properties of interest, and even generate a proposed 'recipe' for how to synthesize the compound. The real-world testing of compounds in biological systems will always be a bottleneck and an important step in the process, but it's exciting to see how much acceleration we can get from computational predictions.

### Wengong Jin

I think AI will be applied to a much broader range of biological applications like structural biology, immunology, gene therapy and drug delivery. Therapeutic development in these areas has been hindered by the enormous time and cost associated with experimental processes. AI-based therapeutic design may become the next-generation technology in these fields.



For example, the success of gene therapy or cancer drugs depends on the efficiency and selectivity of nanoparticles in delivering the drug to desired cell types. We can enhance drug delivery technologies by building neural networks to predict the efficiency and selectivity of nanoparticles and generating new vectors with optimal efficiency and selectivity via generative models.

### Jürgen Bajorath

For the reasons discussed above, I do not anticipate 'revolutionary' AI-driven developments in drug discovery and design over the next years. Provocatively put, making better drugs through AI probably is an elusive goal for the next decade, given that the discovery process is multi-factorial and much too complex and time-consuming for a single technology to be a game-changer.

Instead, incremental advances in early-phase discovery such as in synthesis prediction, targeted compound design, or in vivo drug property predictions are expected and will certainly be helpful. However, for AI/DL to mature in discovery settings, there is an urgent need for more prospective applications (that is, demonstrating what has been accomplished, rather than what could be done). This will primarily depend on the confidence of drug discovery investigators to translate predictions into experiments.

Practical applications in high-profile discovery projects will be essential for establishing AI within the drug discovery spectrum and increasing its acceptance among experimentalists. It is also anticipated that further progress will be made in integrating predictive modeling with robotics in lab automation. Although this might not always require rocket science, the potential impact of such efforts should not be underestimated, especially if they lead to substantial reductions in the workload required for standard procedures in chemical labs, biological screening, or in the scale-up of experiments.

Last but not least, going beyond DL, it will also be very interesting to see a more extensive deployment of other AI methods and tools such as recommender systems that have the potential to impact the practice of drug discovery.

# Meet the Experts

---

Anne Carpenter is the Senior Director of the imaging platform at the Broad Institute (MA, USA) and co-leads a research group developing computational techniques for use across multiple disease areas. Carpenter's research focuses on developing AI methods for biological image analysis. Carpenter acknowledges funding for her laboratory from the National Institutes of Health (R35 GM122547).



Wengong Jin is a Postdoctoral Associate at the Eric and Wendy Schmidt Center at the Broad Institute (MA, USA). Jin's research looks at developing novel machine learning algorithms for a variety of biological applications, most recently using deep learning to identify synergistic drug combinations for treating COVID-19.



Jürgen Bajorath is a Professor and Chair of Life Science Informatics at the University of Bonn (Germany), where his research focuses on developing computational methods for medicinal chemistry, chemoinformatics and chemical biology. Bajorath also studies structure-activity relationships in drug design, Big Data analytics and the use of AI in the life sciences.

