



# HHS Public Access

Author manuscript

*Nat Mach Intell.* Author manuscript; available in PMC 2023 March 07.

Published in final edited form as:

*Nat Mach Intell.* 2022 March ; 4(3): 189–191. doi:10.1038/s42256-022-00465-9.

## Dual Use of Artificial Intelligence-powered Drug Discovery

Fabio Urbina<sup>1</sup>, Filippa Lentzos<sup>2</sup>, Cédric Invernizzi<sup>3</sup>, Sean Ekins<sup>1</sup>

<sup>1</sup>Collaborations Pharmaceuticals, Inc., 840 Main Campus Drive, Lab 3510, Raleigh, NC 27606, USA.

<sup>2</sup>Department of Global Health & Social Medicine, King's College London, United Kingdom.

<sup>3</sup>Spiez Laboratory, Federal Department of Defence, Civil Protection and Sports, Switzerland.

### Abstract

An international security conference explored how artificial intelligence (AI) technologies for drug discovery could be misused for de novo design of biochemical weapons. A thought experiment evolved into a computational proof.

---

The Swiss Federal Institute for NBC-Protection—Spiez Laboratory—is part of the ‘convergence’ conference series<sup>1</sup> set up by the Swiss government to identify developments in chemistry, biology and enabling technologies, which may have implications for the Chemical and Biological Weapons Conventions. Meeting every two years, the conference brings together an international group of scientific and disarmament experts to explore the current state of the art in the chemical and biological fields and their trajectories, to think through potential security implications, and to consider how these implications can most effectively be managed internationally. The meeting convenes for three days of discussion on the possibilities of harm, should the intent be there, from cutting edge chemical and biological technologies. Our drug discovery company received an invitation to contribute a presentation on how AI technologies for drug discovery could be potentially misused.

### Risk of misuse

The thought had never struck us. We were vaguely aware of security concerns around work with pathogens or toxic chemicals, but that did not relate to us; we primarily operate in a virtual setting. Our work is rooted in building machine learning models for therapeutic and toxic targets to better assist in the design of new molecules for drug discovery. We have spent decades using computers and AI to improve human health—not to degrade it. We were naïve in thinking about the potential misuse of our trade, as our aim had always been to avoid molecular features that could interfere with the many different classes of proteins essential to human life. Even our projects on Ebola and neurotoxins, which could have sparked thoughts about the potential negative implications of our machine learning models, had not set our alarm bells ringing.

---

Competing interests

F.U. and S.E. work for Collaborations Pharmaceuticals, Inc. F.L. and C.I. have no conflicts of interest.

Our company—Collaborations Pharmaceuticals, Inc—had recently published computational machine learning models for toxicity prediction in different areas, and, in developing our presentation to the Spiez meeting, we opted to explore how AI could be used to design toxic molecules. It was a thought exercise we had not considered before that ultimately evolved into a computational proof of concept for making biochemical weapons.

## Generation of new toxic molecules

We had previously designed a commercial *de novo* molecule generator which we called MegaSyn<sup>2</sup> which is guided by machine learning model predictions of bioactivity for the purpose of finding new therapeutic inhibitors of targets for human diseases. This generative model normally penalizes predicted toxicity and rewards predicted target activity. We simply proposed to invert this logic using the same approach to design molecules *de novo*, but now guiding the model to reward both toxicity and bioactivity instead. We trained the AI with molecules from a public database using a collection of primarily drug-like molecules (that are synthesizable and likely to be absorbed) and their bioactivities. We opted to score the designed molecules with an organism-specific lethal dose (LD<sub>50</sub>) model<sup>3</sup>, and a specific model using data from the same public database which would ordinarily be used to help derive compounds for treatment of neurological diseases (details of the approach are withheld but were available during the review process). The underlying generative software is built on and similar to other open-source software that is readily available<sup>4</sup>. To narrow the universe of molecules we chose to drive the generative model towards compounds like the nerve agent VX, one of the most toxic chemical warfare agents developed during the 20<sup>th</sup> century—a few salt-sized grains of VX, (6–10 mg)<sup>5</sup>, is sufficient to kill a person. Nerve agents such as Novichoks have also been in the headlines recently<sup>6</sup>.

In less than 6 hours after starting on our in-house server, our model generated forty thousand molecules that scored within our desired threshold. In the process, the AI designed not only VX, but many other known chemical warfare agents that we identified through visual confirmation with structures in public chemistry databases. Many new molecules were also designed that looked equally plausible. These new molecules were predicted to be more toxic based on the predicted LD<sub>50</sub> in comparison to publicly known chemical warfare agents (Figure 1). This was unexpected as the datasets we used for training the AI did not include these nerve agents. The virtual molecules even occupied a region of molecular property space that was entirely separate to the many thousands of molecules in the organism-specific LD<sub>50</sub> model, which is mainly made up of pesticides, environmental toxins, and drugs (Figure 1). By inverting the use of our machine learning models, we had transformed our innocuous generative model from a helpful tool of medicine to a generator of likely deadly molecules.

Our toxicity models were originally created for use in avoiding toxicity, enabling us to better virtually screen molecules (for pharmaceutical and consumer product applications) before ultimately confirming their toxicity through *in vitro* testing. The inverse, however, has always been true: the better we can predict toxicity, the better we can steer our generative model to design new molecules in a region of chemical space populated by predominantly lethal molecules. We did not assess the virtual molecules for synthesizability or explore how

to make them with retrosynthesis software. Both of these processes have readily available commercial and open-source software, which can be easily plugged into the *de novo* design process of new molecules<sup>7</sup>. We also did not physically synthesize any of the molecules either, but with a global array of hundreds of commercial companies offering chemical synthesis, it is not necessarily too big of a step, which is poorly regulated with few if any checks to prevent synthesis of new extremely toxic agents that could potentially be used as chemical weapons. Importantly, we had a human-in-the-loop with a firm moral and ethical ‘don’t-go-there’ voice to intervene. But what if the human was removed or replaced with a bad actor? With current breakthroughs and research into autonomous synthesis<sup>8</sup>, a complete design-make-test cycle applicable to making not only drugs, but toxins, is within reach. Our proof-of-concept highlights how a non-human autonomous creator of a deadly chemical weapon is entirely feasible.

## A wake-up call

Without being overly alarmist, this should serve as a wake-up call for our colleagues in the ‘AI in drug discovery’ community. While some domain expertise in chemistry or toxicology is still required to generate toxic substances or biological agents that can cause significant harm, when these fields intersect with machine learning models, where all you need is the ability to code and to understand the output of the models themselves, they dramatically lower technical thresholds. Open source machine learning software is the primary route for learning and creating new models like ours, and toxicity datasets<sup>10</sup> that provide a baseline model for predictions for a range of targets related to human health are readily available.

Our proof of concept was focused on VX-like compounds, but it is equally applicable to other toxic small molecules with similar or different mechanisms with minimal adjustments to our protocol. Retrosynthesis software tools are also improving in parallel, allowing new synthesis routes to be investigated for known and unknown molecules. It is therefore entirely possible that novel routes can be predicted for chemical warfare agents, circumventing national and international lists of watched or controlled precursor chemicals for known synthesis routes.

The reality is that this is not science fiction. We are but one very small company in a universe of many hundreds of companies using AI software for drug discovery and *de novo* design. How many of them have even considered repurposing, or misuse, possibilities? Most will work on small molecules and many of the companies are very well funded and likely using the global chemistry network to make their AI designed molecules. How many people are familiar with the know-how to find the pockets of chemical space that can be filled with molecules predicted to be orders of magnitude more toxic than VX? We do not currently have answers to these questions. There has not previously been significant discussion in the scientific community about this dual use concern of AI used for *de novo* molecule design, at least not publicly. Discussion of societal impact of AI has principally focused on aspects like safety, privacy, discrimination and potential criminal misuse<sup>10</sup>, but not national and international security. When we think of drug discovery, we normally do not consider technology misuse potential. We are not trained to consider it, and it is not even required for machine learning research, but we can now share our experience with other

companies and individuals. AI generative machine learning tools are equally applicable to larger molecules (peptides, macrolactones etc.) and to other industries like consumer products and agrochemicals that also have interests in designing and making new molecules with specific physicochemical and biological properties. This greatly increases the breadth of the potential audience that should be paying attention to these concerns.

For us, the genie is out of the medicine bottle when it comes to repurposing our machine learning. We must now ask: what are the implications? Our own commercial tools as well as open-source software tools and many datasets that populate public databases are available with no oversight. If the threat of harm, or actual harm, occurs with ties back to machine learning, what impact will this have on how this technology is perceived? Will hype in the press on AI-designed drugs suddenly flip to AI-designed toxins, public shaming, and decreased investment in these technologies? As a field, we should open a conversation on this topic. The reputational risk is substantial; it only takes one bad apple that takes what we have vaguely described to the next logical step, or an adversarial state looking for a technological edge. How do we prevent this? Can we lock away all the tools and throw away the key? Do we monitor software downloads or restrict sales to certain groups? We could follow the example of machine-learning models like GPT-3<sup>11</sup> which was initially waitlist restricted to prevent abuse and has an API for public usage. Even today, without a waitlist, GPT-3 has safeguards in place to prevent abuse, Content Guidelines, a free content filter and monitoring of applications that use GPT-3 for abuse. We know of no recent toxicity or target model publications that discuss these concerns of dual use, similarly. As responsible scientists, we need to ensure that misuse of AI is prevented, and that the tools and models we develop are only used for good.

By going as close as we dared, we have still crossed a grey moral boundary, demonstrating that designing virtual potential toxic molecules is possible without much effort, time or computational resources. We can easily erase the thousands of molecules we created, but we cannot delete the knowledge of how to recreate them.

## The broader impacts on society

There is a need for discussions across traditional boundaries and multiple disciplines to allow for a fresh look at AI for *de novo* design and related technologies from different perspectives and with a wide variety of mindsets. Here, we give some recommendations which we believe will reduce potential dual-use concerns for AI in drug discovery. Scientific conferences, like the Society of Toxicology and American Chemical Society, for example should actively foster a dialogue among experts from industry, academia and policy making on the implications of our computational tools. There has been recent discussion in this journal regarding requirements for broader impact statements from authors submitting to conferences, institutional review boards and funding bodies as well as addressing potential challenges<sup>12</sup>. Making increased visibility a continuous effort and a key priority would greatly assist in raising awareness about potential dual use aspects of cutting-edge technologies and would generate the outreach necessary to have everyone active in our field engage in responsible science. We can take inspiration from examples such as The Hague Ethical Guidelines<sup>13</sup>, which promote a culture of responsible conduct in the chemical

sciences and guard against the misuse of chemistry, in order to have AI-focused drug discovery, pharmaceutical, and possibly other companies agree to a code of conduct to train employees, secure their technology and prevent access and potential misuse. The use of a public-facing API for models with code and data available upon request would greatly enhance the security and control over how published models are utilized without adding much hindrance to accessibility. While MegaSyn is a commercial product and thus we have control over who has access to it, going forward, we will implement restrictions or an API for any forward-facing models. A reporting structure or hotline to authorities, if there is a lapse or if we become aware of anyone working on developing toxic molecules for non-therapeutic uses, may also be valuable. Finally, universities should redouble their efforts in ethical training of science students and broaden the scope to other disciplines, particularly computing students, so that they are aware of the potential misuse of AI from an early stage of their career as well as understand the potential for broader impact<sup>12</sup>. We hope that by raising awareness of this technology we will have gone some way to demonstrating that while AI can have important applications for healthcare and other industries, we should also remain diligent against the potential for dual use, in the same way that we would with physical resources such as molecules or biologics.

## Acknowledgments

We are grateful to the organizers and participants of the Spiez Convergence conference 2021 for their feedback and questions.

Cédric Invernizzi contributed to this article in his personal capacity. The views expressed in this article are those of the authors only and do not necessarily represent the position or opinion of Spiez Laboratory or the Swiss Government.

## Funding

We kindly acknowledge NIH funding from R44GM122196-02A1 from NIGMS and 1R43ES031038-01 and 1R43ES033855-01 from NIEHS for our machine learning software development and applications. "Research reported in this publication was supported by the National Institute of Environmental Health Sciences of the National Institutes of Health under Award Number R43ES031038 and 1R43ES033855-01. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health."

## References

1. <<https://www.spiezlab.admin.ch/en/home/meta/refconvergence.html>> (2021).
2. Urbina F, Lowden CT, Culberson JC & Ekins S <[10.33774/chemrxiv-2021-nlwvs](https://doi.org/10.33774/chemrxiv-2021-nlwvs)> (2021).
3. Mansouri K et al. Environ Health Perspect 129, 79001, (2021). [PubMed: 34242083]
4. Blaschke T et al. J Chem Inf Model 60, 5918–5922, (2020). [PubMed: 33118816]
5. Anon <<https://www.ncbi.nlm.nih.gov/books/NBK233724/>> (1997).
6. Aroniadou-Anderjaska V, Apland JP, Figueiredo TH, De Araujo Furtado M & Braga MF Neuropharmacology 181, 108298, (2020). [PubMed: 32898558]
7. Genheden S et al. J Cheminform 12, 70, (2020). [PubMed: 33292482]
8. Coley CW et al. Science 365, eaax1566, (2019). [PubMed: 31395756]
9. Dix DJ et al. Toxicol Sci 95, 5–12, (2007). [PubMed: 16963515]
10. Hutson M <<https://www.newyorker.com/tech/annals-of-technology/who-should-stop-unethical-ai>> (2021).
11. Brown TB et al. <<https://arxiv.org/abs/2005.14165>> (2020).
12. Prunkl CEA et al. Nature Machine Intelligence 3, 104–110, (2021).

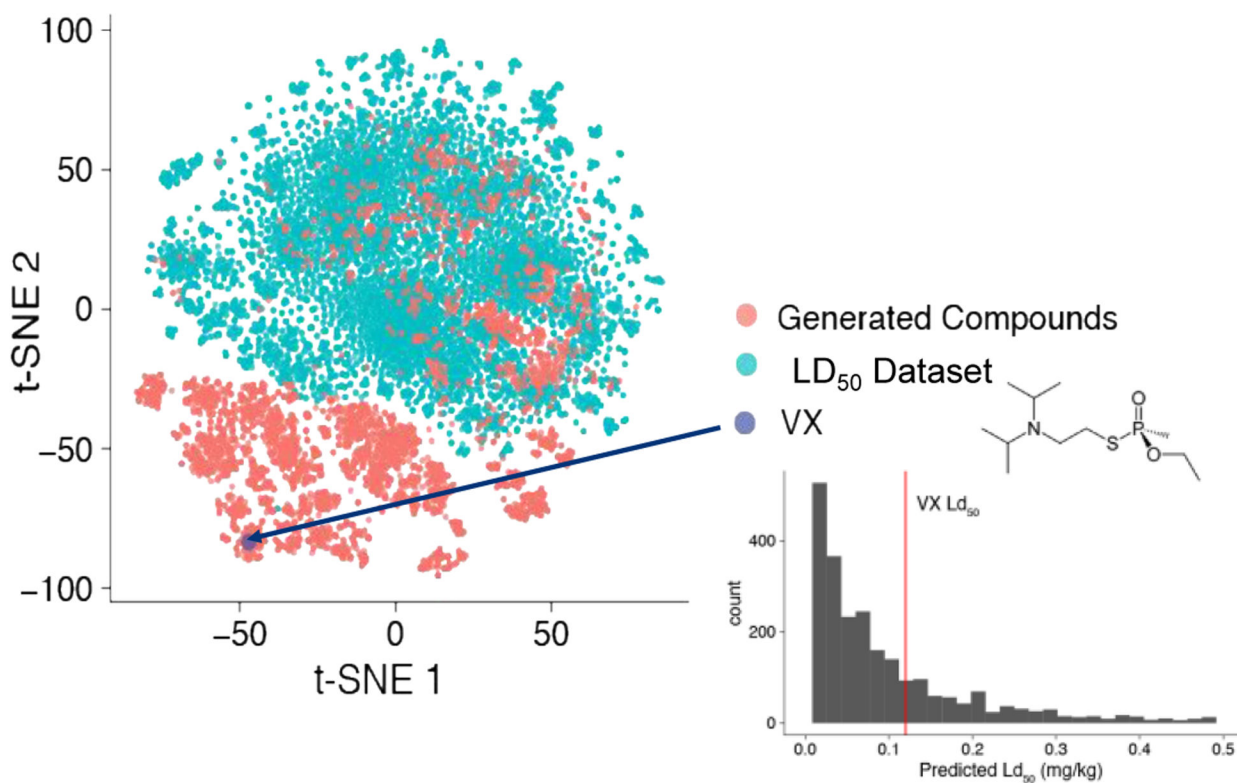
13. <<https://www.opcw.org/hague-ethical-guidelines>> (2021).

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript



**Figure 1.**

A t-SNE plot visualization of the LD<sub>50</sub> dataset (cyan) and top 2000 MegaSyn AI-generated and predicted toxic molecules (salmon) illustrating VX (purple and 2D structure). Many of the molecules generated are predicted as more toxic *in vivo* in the animal model than VX (histogram showing cutoff for VX LD<sub>50</sub>).