

# Assessing Productivity of Pharmaceutical AI: Any Results Beyond Hype?

March 18, 2019 by Andrii Buvailo

## A background context -- opportunities and challenges

Current widespread interest towards artificial intelligence (AI) and its numerous research and commercial successes was largely catalyzed by several landmark breakthroughs in 2012, when researchers at the University of Toronto achieved unprecedented improvement in the image classification challenge ImageNet, using their deep neural network “AlexNet” running on graphics processing units (GPUs), and when that same year Google’s deep neural network managed to identify a cat from millions of unlabeled Youtube videos, representing a conceptual step in unsupervised machine learning.

Nowadays, AI is commercially used in countless applications across various industries, ranging from surveillance, finance, marketing, automation, robotics, driverless cars etc. Pharmaceutical industry appeared to be a relative late adopter of AI tech, but the progress here is catching up rapidly now, for example, there were more research publications on the “AI for drug discovery” in 2018 alone, than in all previous years combined.

**(Read also: “2018: AI Is Surging In Drug Discovery Market”)**

However, there is a major bottleneck in the application of deep learning (the current “workhorse” of modern AI) -- the need for large amounts of clean and properly linked data. The famous dataset ImageNet, which led to the AlexNet’s success, included 15 million labeled high-resolution images in over 22,000 categories, which represented a high quality big data set. While in some cases obtaining clean data in large volumes is a manageable task, like in the case of speech recognition, in other cases data is scarce and disperse across multiple poorly linked sources -- which is the case in life sciences industry. Surprisingly, there is a tendency to overestimate the amount of available quality data in pharmaceutical industry. Not only a lot of research data in drug development is in many cases poorly validated, but even properly validated data in a possession of numerous biopharma companies might not be easily available due to an inherent and strict code of secrecy, defined by severe competition among drug makers. In this situation, combining enough drug development data from multiple organizations and projects to really take advantage of AI/ML technologies on a global scale becomes a considerable challenge for the industry.

While the lack of quality data for training machine learning models might be, in some instances, overcome by state-of-the-art strategies, like transfer learning, synthetic data, or approaches to train on small data sets, the “prime-time” has not yet come for them in the pharmaceutical industry.

**(Read also: “Democratizing Artificial Intelligence For Pharmaceutical Research”)**

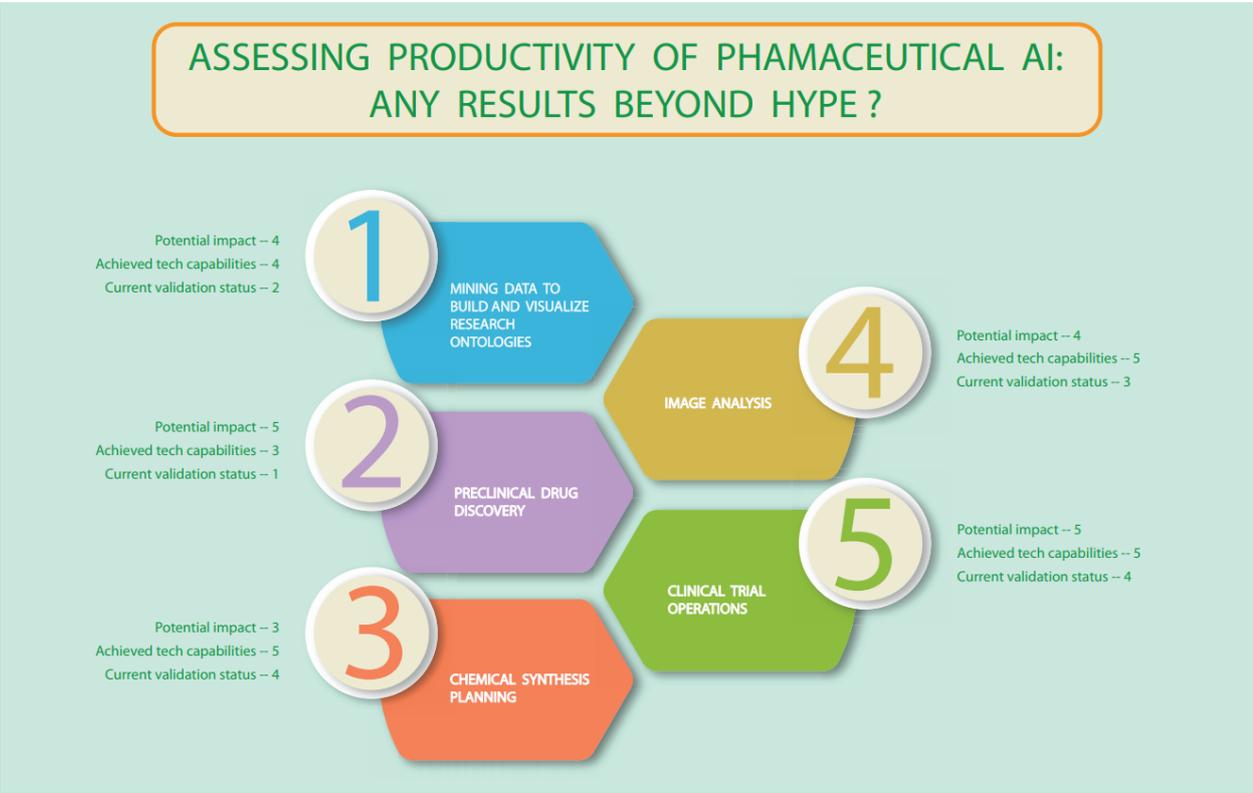
Another bottleneck -- validation of any findings offered by AI after performing analysis (especially, in unsupervised learning mode) on highly domain-specific and feature-rich datasets -- e.g. those derived from “omics” studies in biomedical research programs. If we talk about identifying cats in millions of unannotated videos -- that is fine, as we all know cats, and we can appreciate results right away. But when AI finds intricate data patterns or correlations in feature-rich biomedical data, it might not be obvious what we will be looking at. So far, AI can find both meaningful patterns, and nonsense patterns -- just as well.

This issue is reflected in an interview with Dr. Alex Zhavoronkov, one of the leading experts in the area of pharmaceutical AI: “When applying the deep learning techniques to images and videos validation is almost instantaneous and often you can see what to fix. This rapid pace of validation allowed for the driverless Tesla and for the deep-learned Google translate. But when dealing with molecules or biomarkers, it takes a very long time.”

Notwithstanding the above bottlenecks, AI is already delivering value in a wide range of pharmaceutical tasks, albeit not in all cases as rapidly as numerous media outlets tend to predict. Let’s review and assess some practical outputs achieved so far with the adoption of AI.

(For every use case, I assessed three parameters below, applying points from 1 to 5:

- a)** state-of-the-art capabilities of AI technology in a particular use case,
- b)** validation status in practice (public disclosures, FDA-approvals, user references etc), and
- c)** potential business impact of that particular use case on the pharmaceutical industry in general.



(Infographic design by Evgeniy Gumenyuk)

# 1. Mining data to build and visualize research ontologies

**Achieved tech capabilities -- 4**

**Current validation status -- 2**

**Potential impact -- 4**

If we are to talk about one use case where AI can apparently be valuable for life scientists right now, a good start is in mining data from research literature, patents, and other unstructured (and structured) sources to be able to derive research insights augmenting thought process. After all, Google, Facebook,

and many other companies have already cracked this use case and integrated AI in their search services, showing striking accuracy for a broad range of search and contextualization applications.

Among companies, specifically focused on applying AI for mining life science data from diverse sources and using it to provide actionable insights, there are several notable examples, including a Google-backed BenchSci (an AI-powered antibody search service), and Bioz (the world's first AI-driven search engine for life sciences). Another example is a rapidly growing BenevolentAI, a British biotech with several own drug candidates in development. While BenevolentAI is focused on a wide range of drug discovery tasks, one of its core AI strength is in the ability to “mine new knowledge from vast quantities of biomedical data”, which apparently contributed to a company's rapid valuation growth. A strong case of using vast structured and unstructured public and private data sources for AI-driven hypotheses generation is presented by a Baltimore-based Insilico Medicine, a leading AI-company in the area of drug discovery, biomarker development, and aging research.

The use of AI for mining data and building research ontologies is an important and immediate game changer for the pharmaceutical industry for several reasons. First, the area of natural language processing (NLP), relevant here, is among the most advanced in AI field, so the technology is already powerful enough to have practical impact right now. Second, there is plenty of peer-reviewed (hence, relatively validated) data out there to “play” with, so machines can learn efficiently here and up to high levels of accuracy. Third, the validation of results is relatively straightforward (as with cats in the Google experiment), the model can be updated quickly. Finally, every pharmaceutical (and not only) researcher in the world is struggling with exponentially growing number of publications, patents and other data to deal with in their daily work, so having some intelligent and automated tools to assist in this process is becoming essential to increase R&D efficiency.

## 2.Preclinical drug discovery

**Achieved tech capabilities -- 3**

**Current validation status -- 1**

**Potential impact -- 5**

While mining data and generating hypotheses is a very general use case, augmenting almost every stage of drug discovery workflow, and in a very wide range of contexts, there are more specific and more potentially disruptive use cases of AI application, like those for, say, deconvoluting novel promising targets from functional genomics, high content screening or other omics programs; or generating novel drug-like molecules via de-novo molecular design; or accurately predicting safety of drug candidates at early stages of development. This kind of use cases can have transformative impact on any pharmaceutical company, potentially reducing R&D failure rates, time and cost of drug discovery.

While applying AI for preclinical drug discovery is a world of opportunities, we stumble upon the above discussed issues with insufficient, poorly linked, or simply locked by secrecy research data, and the difficulty in validating AI-derived predictions.

Indeed, for the AI to be able to derive a promising new drug molecule from scratch it needs to have a quite comprehensively constructed model of all biological processes and pathways, including disease pathways, and rely heavily on properly validated experimental data from real-world assays and tests. Otherwise, any predictions are not very far from a wild guess.

There are only a handful of AI-driven startups out there which seem to have approached quite close to having integrated “end-to-end” AI-driven drug discovery platforms, starting from hypothesis generation and all the way to optimizing leads and predicting their properties. Three notable examples include Exscientia, a Dundee-based company, claiming to have become a “full stack AI drug discovery company” (and having an impressive portfolio of research deals with big pharma, results TBD), the Insilico Medicine, claiming to have an “end-to-end” AI engine for drug discovery, which recently licensed several small molecules for further development to Juvenescence Ltd (official results TBD), and Healx -- an AI-driven biotech, which applied its full-cycle research platform to come up with a clinical-stage drug candidate for rare diseases.

**(UPDATE:** Exscientia just announced that it has delivered the first AI-derived drug candidate for their collaboration with GSK).

Some other companies probably belonging to this “AI generalist” club include Cyclica, Elucidata, and TwoXAR -- all building integrated AI platforms for a wide range of drug discovery tasks.

Most other companies of the “AI for drug discovery” cohort appear to be more specific about the use cases they focus on. Such examples include Atomwise (structure based virtual screening); Intomics (target identification and validation from omics data); Envisagenics (focusing on RNA), ReviveMed (focusing on metabolomics) etc.

**(Refer to “143 Startups Applying Artificial Intelligence (AI) for Drug Discovery and Biomarker Development” for more information)**

Several of the “AI-driven” biotechs have already progressed to clinical trials with drug candidates, supposedly derived from direct application of ML/AI tools. Those include Berg Health, BenevolentAI, BioXcel, Healx, Lantern Pharmaceuticals, and Recursion Pharmaceuticals. But it is hard to say, to which extent methods like deep learning played roles in those successes. Was it a crucial technical enabler in every case? How would companies have done in relative terms without ML/AI platforms with legacy analytics tools? Remains TBD. But it seems impressive that Healx managed to come up with a Phase 2a clinical trial candidate within 15 months with just \$100,000 budget. With that, I do believe AI will eventually have its “prime-time” in the preclinical drug discovery.

### 3. Chemical synthesis planning

**Achieved tech capabilities -- 5**

**Current validation status -- 4**

**Potential impact -- 3**

Almost all small molecule drug discovery is heavily dependant on efficient chemical synthesis, being in the core of most hit identification, prioritization, and lead optimization programs. The application of various machine learning methods is showing pretty amazing results here -- the state of the art AI-driven tools can now outperform best synthetic chemists in planning complex organic synthesis.

One of the most known AI-driven commercial software for retrosynthesis planning is Chematica, acquired by Merck in 2017. Not only it can help chemists find cheaper and more efficient synthetic routes to complex molecules, but also, it can help avoid patented synthetic routes.

Another notable example is a deep neural network-based tool developed by Marwin Segler, an organic chemist and artificial-intelligence researcher at the University of Münster in Germany, which showed impressive results in planning organyc synthesis.

Here is a nice review, published in Nature, about latest achievement of AI in organic chemical synthesis, and some references to results validation.

Overall, the striking success of AI-based tools in chemical synthesis planning are due to the abundance of relatively clean and quality data to train models, an understandable feature engineering process, and a relatively straightforward result validation. The prime time of AI is arriving for this use case.

## 4. Image analysis

**Achieved tech capabilities -- 5**

**Current validation status -- 3**

**Potential impact -- 4**

Medical image analysis is a “holy grail” of AI application in healthcare (primarily diagnostics), but in pharmaceutical research it can also serve well in a number of use cases, for instance, in multimodal high throughput image analysis for quantification of structural and functional evolution of cells, tissues and organs.

In this regard, one of the leading AI-driven companies applying cell image analysis for discovering novel drugs is Recursion Pharmaceuticals, a biotech company headquartered in Salt Lake City. The company developed an automatic massively parallel AI-based platform that performs large-scale analytics of

cellular phenotypes to develop computational “fingerprints” (Phenoprints™) of a wide variety of biological perturbations. Needless to say, Recursion has own pipeline of drug candidates with two of them in the Phase 1 clinical trials.

An important insight about applying AI in image analysis came from Novartis CEO Vas Narasimhan in this interview podcast -- about a massive project to digitize all of the company’s pathology images by partnering with PathAI, and to expand this process for other categories of images as well. He notes that this is a work in progress, which value is still TBD.

## 5.Clinical trial operations

**Achieved tech capabilities -- 5**

**Current validation status -- 4**

**Potential impact -- 5**

Running efficient clinical trials is the central task of any pharmaceutical company, a major milestone towards overall business success. Currently, AI is being applied for patient recruitment (for instance, Antidote, Deep 6 AI and Mendel.ai), clinical trial design (Trials.AI and BullFrog AI), and clinical trial optimization (Brite Health).

There are reports confirming value of the AI application in clinical trials, for example, in this review a comparison of AI-powered vs standard methods was made for three oncology trials, and it was found, for example, that for each trial that enrolled, the use of Mendel.ai resulted in a 24% to 50% increase over standard practices in the number of patients correctly identified as potentially eligible. The conclusion was that AI can assure measurable improvements over standard methods in several use cases, including patient pre-screening process, feasibility, site selection, and trial selection.

Mayo Clinic managed to increase patient enrollment by 80% using IBM Watson for Clinical Trial Matching platform (amidst some public criticism of this particular AI tool).

In this interview podcast Novartis' CEO Vas Narasimhan explains that the company organized a high-tech control center for monitoring and planning all of their clinical trials in the world, where AI is predicting patient enrollment statistics, and assessing possible quality issues -- all this based on ten years of historical data for previous clinical trials, as well as new data coming from 400-500 clinical trials conducted every year, available to train the AI models to become valuable.

So again, the abundance of quality data and a relatively quick validation cycle allows for an immediate success of AI application in assisting clinical trial operations.

## Conclusions

AI (primarily deep neural nets) is a transformative technology in many ways, which is validated by numerous practical applications outside pharmaceutical industry.

AI is increasingly being adopted by pharmaceutical industry and it is already showing measurable value for a number of use cases, including search and ontology building, augmenting mundane research tasks, controlling clinical trials etc. While for the more domain-specific tasks, like in preclinical drug discovery, AI is still in a validation stage as a game-changing solution, might take 2-3 years for the FDA-quality results to emerge.

The fact that AI works staggeringly well for one use case does not necessarily mean it will work the same for another use case. Applicability depends (at least) on the abundance of clean and quality data to train models, and the ability to validate results efficiently.

As Amara's law states: **"We tend to overestimate the effect of a technology in the short run and underestimate the effect in the long run"**. Very true for the AI in pharmaceutical research.