



HHS Public Access

Author manuscript

Pharm Res. Author manuscript; available in PMC 2017 November 01.

Published in final edited form as:

Pharm Res. 2016 November ; 33(11): 2594–2603. doi:10.1007/s11095-016-2029-7.

The Next Era: Deep Learning in Pharmaceutical Research

Sean Ekins*

Collaborations Pharmaceuticals, Inc., 5616 Hilltop Needmore Road, Fuquay-Varina, NC 27526, USA; Collaborative Drug Discovery, 1633 Bayshore Highway, Suite 342, Burlingame, CA 94010, USA

Abstract

Over the past decade we have witnessed the increasing sophistication of machine learning algorithms applied in daily use from internet searches, voice recognition, social network software to machine vision software in cameras, phones, robots and self-driving cars. Pharmaceutical research has also seen its fair share of machine learning developments. For example, applying such methods to mine the growing datasets that are created in drug discovery not only enables us to learn from the past but to predict a molecule's properties and behavior in future. The latest machine learning algorithm garnering significant attention is deep learning, which is an artificial neural network with multiple hidden layers. Publications over the last 3 years suggest that this algorithm may have advantages over previous machine learning methods and offer a slight but discernable edge in predictive performance. The time has come for a balanced review of this technique but also to apply machine learning methods such as deep learning across a wider array of endpoints relevant to pharmaceutical research for which the datasets are growing such as physicochemical property prediction, formulation prediction, absorption, distribution, metabolism, excretion and toxicity (ADME/Tox), target prediction and skin permeation, etc. We also show that there are many potential applications of deep learning beyond cheminformatics. It will be important to perform prospective testing (which has been carried out rarely to date) in order to convince skeptics that there will be benefits from investing in this technique.

Keywords

Artificial intelligence; Deep Learning; Drug Discovery; Machine learning; Pharmaceutics

Introduction

We have previously suggested that cheminformatics should look to other industries that use high performance computing approaches for inspiration (1). Six years on what is surprising is that we may not have to look to industries but instead we already have used the algorithms which we should be considering for cheminformatics (and other areas) in our everyday transactions, work or social-life online. Every time we use our smartphones with voice recognition software like Siri or read the news on them, make a purchase on the internet via

*To whom correspondence should be addressed. ekinssean@yahoo.com.

Conflict of interest

SE is founder and owner of Collaborations Pharmaceuticals, Inc. he was a consultant for Collaborative Drug Discovery, Inc.

Amazon or use our social network software, we take for granted that we are confronted with suggestions of other products we might like to buy or friends to connect too. Large companies such as Baidu, Google, Facebook etc. all use deep learning in facial recognition algorithms alone. We also live in an age in which self-driving cars and robot assistants are emerging rapidly after decades of research. Therefore we are literally surrounded by artificial intelligence software that use machine learning (2) that can in many ways now predict our needs before we know what they are. Perhaps for the first time we may also be directly seeing a preview of how such software tools could assist us in healthcare related research and development. We just did not know it until now.

In the area of pharmaceutical research and development and specifically that of cheminformatics there are many machine learning methods such as support vector machines (SVM), k-Nearest Neighbors, Naïve Bayesian, Decision Trees etc. (3) which have seen increasing use as our datasets have grown to become 'big data' (4-7). These methods are equally applicable in other areas and come with their own pros and cons (3) that can be used for binary classification, multiple classes or continuous data. The application of different computational approaches and machine learning algorithms to problems tends to follow the growth of datasets (8). As pharmaceutical datasets started out quite small, the methods initially used would focus on local models like pharmacophores and quantitative structure activity relationships (QSAR). In more recent years the biological data amassed from high throughput screening and high content screens has called for different tools to be used that can account for some of the issues with this bigger data (6). The power of computer processing has also increased so that more complex non-linear problems can be solved in real time with relatively inexpensive compute resources. Many of these resulting machine learning models can also be implemented on a mobile phone (9, 10). In recent years, there has been increasing use of one approach called deep learning (DL), (which builds on many years of artificial neural network research)(11), that has shown powerful advantages in learning from images and languages (12). This may represent the next era of cheminformatics and pharmaceutical research in general that is focused on mining the heterogeneous big data which is accumulating using more sophisticated algorithms such as DL.

Delving into deep learning

Standard artificial neural networks (ANN) approaches use an input layer, hidden layer and output layer where each connection has a weight and these vary during training in order to connect input to output data. This method has been used widely but suffers from overfitting of data, and a poor ability to generalize with an external dataset (3), although more recent versions such as Bayesian regularized artificial neural networks are more difficult to overtrain (13). DL or deep neural networks (3) in many ways is similar to ANN in that it mimics how the brain works and takes information in an input layer but unlike ANN has many hidden layers (14) to combine signals with different weights, passes the results successively deeper in the network until an output layer (Figure 1). The DL model is trained with a dataset by adjusting the weights to give the response expected for a certain input (e.g. if a compound is active or inactive or the level of activity/inactivity). The ability to have multiple learnable stages makes this approach more useful for tackling more complex

problems. Deep learning can be used for unsupervised learning and appears to work well with noisy data. However it still suffers from the potential to over fit data, the black box problem, as well as higher computational cost than ANN or other methods (15). There has been relatively limited application of DL to pharmaceutical problems to date and very few in the area of cheminformatics compared with other machine learning methods (11). DL tools are available in popular open source statistical software such as R (16). In addition there is TensorFlow (17), Deeplearning4j (18) while Facebook made their deep learning software (Torch) open source (19, 20) followed a year later by Microsoft (CNTK) (21). Some of these methods have been summarized in a recent review (22). It should be noted that these deep learning toolkits are likely far from 'plug and play' type software tools for the average scientist which they can input their molecules and data to train a model (or for that matter any training or test datasets) and then generate predictions. It is likely that expertise in using these software toolkits is needed as well as integrating with molecular descriptor software. It is more likely that a specialized programmer / statistician / cheminformatician will be needed with knowledge of the software tools in order to generate the models which can then be made available for others to use. Existing cheminformatics and other software companies could facilitate making deep learning more accessible to non-expert users by developing accessible fully integrated tools which they can use for any dataset.

Applications of deep learning in bioinformatics

DL has seen a rapid increase in the number of publications associated with bioinformatics (23) and computational biology (22) and has been used in diverse applications (15) such as protein disorder prediction (24). This resulted in a fast approach which was comparable to other machine learning approaches based on area under the curve (AUC) and recall statistics for test datasets. While it did not out-perform the other disorder methods it had the advantage of being fast (24) (Table 1). DL has also been used to refine docked protein complexes (25) based on 35,000 unbound docking complexes generated by RosettaDock and tested on 25 docking complexes not in the training set. Although this model was not compared to other methods it resulted in RMSD of 1.40Å indicating accurate predictions (25). DL has also been used to model structural features of RNA-binding protein targets using CLIP-seq datasets and testing on 24 datasets using the area under the receiver operator characteristic to evaluate the performance which was found to be close to the state of the art (Table 1)(26). DL has been applied and compared with other methods for mechanism of action prediction from high content image analysis data (27) and it was found to be superior to SVM with (87.62 vs 20.95% accuracy). Although it should be noted that the processing time was ten times longer for DL which is likely a major limitation for learning versus other methods (27). The preprocessing time can likely be offset by parallelization of calculation processes.

Pharmaceutical applications of deep learning

One of the earliest applications of DL to a pharmaceutically relevant problem was to predict aqueous solubility using four published datasets and was shown to compare favorably to other machine learning methods using 10 fold cross validation (28) (Table 1). DL has been used to predict the site of epoxidation in molecules with average AUC > 94.0% for cross

validation, although this method was not compared to additional machine learning methods (29) or used for prospective prediction. DL has been put to use with gene expression data to learn from drugs and therapeutic categories using pathway level or landmark gene level as data reduction methods. In both cases after 10 fold cross validation deep neural networks surpassed SVM used internal testing, suggesting this as a drug repurposing approach (Table 1)(30). In the area of drug formulation, predicting drug release from poly-lactide-co-glycolide (PLGA) microspheres showed deep learning to be comparable to random forest, single tree and genetic algorithms after 10 fold cross validation (31). At Merck, deep neural networks have been compared to random forests for use with large QSAR datasets and outperformed random forests for 11 out of 15 datasets (used in a Kaggle crowdsourcing competition) and 13 of 15 datasets in a second evaluation using time-split test sets (32). This utilization by a major pharmaceutical company suggests there is serious interest in the approach and it is likely other companies may have already performed similar evaluations. More frequently DL is applied to a single dataset such as drug induced liver injury. In this case multiple training and testing sets were used and comparison with normal neural networks was performed showing slight improvement with DL (33). In addition, one model was tested with 6 datasets as a form of external validation. The accuracy of DL with DILI data was around 60% for these test sets which is comparable to what has been seen elsewhere using other algorithms (34). Based on a non-exhaustive assessment of several different end points relevant to pharmaceutical research, while it appears that most have seen utilization of Bayesian or SVM approaches to develop predictive models, few have so far utilized DL (Table 2). Recent examples of computational models appearing in this journal alone over the past 18 months include: modeling thermodynamic proxies (35), predicting mouse liver microsomal stability (36), predicting autooxidation (37), drug solubility in human intestinal fluid (38), site of metabolism prediction in CYP2C9 (39), human skin permeability prediction (40, 41), blood brain barrier penetration modeling (42), predicting clearance mechanism (41) and skin concentration due to dermal exposure (43). Many of these datasets could likely utilize and benefit from DL and it would be of interest to see for how many an improvement in predictions could be obtained.

The future of Deep Learning

While there has been recent exhaustive analysis of artificial intelligence and its impact on jobs, ethical considerations and geopolitical impact (44) there have been very few discussions of the potential for using DL in pharmaceutical research (14, 15). Based on the results obtained to date which admittedly have focused on internal validation with little prospective testing as seen with other machine learning methods (Table 1)(45, 46), DL appears promising and will likely see greater application in the years ahead. Perhaps the largest example of validation of DL models alongside other machine learning approaches is in the case of the Tox21 Challenge. DL with multitask learning (47) slightly outperformed the closest consensus artificial neural network method (48) across nuclear receptor and stress response datasets (Table 1).

It is yet to be seen if DL could facilitate the ultimate robot scientist (49), as we see application to different datasets, it or other machine learning methods may become an invisible research assistant, a tool that we take for granted to perform the predictions we

need before performing the experiments. This may come with challenges such as how much power do we provide to the software to make the decisions for us. Which raises the question of whether the experiments need human involvement at all? The rapid development of DL outside science suggests it is far from a new fad, and the impact is already being felt in numerous areas from fraud detection to internet search engines. So how long before DL is widespread in pharmaceutical research (14) and what can we expect? It is possible that DL could be the source of more predictive models but hurdles remain on the implementation and accessibility of models. What is clearly needed is software that is tightly integrated with the data to be modeled. This data would most frequently reside in private or public databases and could represent many different endpoints both quantitative and qualitative (Figure 2). Therefore any efforts to bring the molecules, sources of data and DL algorithms together would greatly streamline model generation and make it more accessible to other scientists. But as with other computational modeling approaches we may also want to consider the applicability domain (50) and various factors such as the quality of the underlying data (51, 52) which may determine the utility and relevance of a DL model for making a prediction (53).

A major concern would be how could DL shape research and the future of science and biomedical research in particular? It is possible that DL or any machine learning method might be able to assist by increasing the efficiency of research and perhaps rule out likely less successful avenues of research. This is especially important in areas where research funding is tight like rare (54) and neglected diseases (55). Any advantage that DL could provide would be welcome in these and other resource constrained areas as a side effect. Clearly we should be educating the next generation of pharmaceutical researchers to use a wide array of machine learning approaches as well as assessing the likely impact and application of DL. Developing scientists that can generate predictions *in silico* and test them *in vitro* or *in vivo* would also be welcomed. It may be only a short time before we have vice presidents of machine learning or DL in pharmaceutical and biotechnology companies. There are of course still many criticisms of these black box machine learning approaches but it is probably now accepted that with greater accuracy in prediction will come limitations in transparency.

While we have only just seen the beginning of the era of DL, we should be prepared for how it will be used and its potential impact on a wide array of potential pharmaceutically relevant endpoints (Table 2). Already comparisons of DL with additional machine learning algorithms have shown that it frequently improves upon the state of the art using predominantly cross validation as the form of evaluation. We will likely see improvements in DL and maybe even alternative approaches that are superior by combining with other machine learning or other methods or data as consensus approaches. This in particular is an area we have yet to see developed for pharmaceutical applications. If we are to imagine machine learning models being ubiquitous in the pharmaceutical industry, DL may facilitate that. We predict the near future will likely see an increase in studies published in this and other journals applying DL and comparing it prospectively and retrospectively to other machine learning methods for predicting various molecule properties important for pharmaceutical research.

While at the time of writing there are over 100 DL startup companies globally, few are focused on pharmaceutical applications (56, 57). We anticipate that this will be an active area as the DL connection will be one way to attract technology investors who would normally steer clear of drug discovery and pharmaceuticals as an investment. It is likely that researchers in academia and industry could immediately apply DL if it was more accessible in software and they could plug their data into it. In the interim it may require some collaboration with those more experienced data scientists using R and the various available DL toolkits. There are several areas that could perhaps see an immediate benefit, for example in the areas of predicting metabolism, interactions with P450s and transporters as well as other ADME/Tox properties there has been an increasing number of large datasets and computational machine learning models published (Table 2) some of which are non-proprietary. Big datasets like melting point with over 300,000 compounds (58), DMSO solubility with over 163,000 compounds (59), all of the ChEMBL (5) datasets and of course data in PubChem (60), would represent “Big Data” (14) that could be used with DL. The data from PubChem and to a lesser extent ChEMBL, may need curation and organization prior to use in DL models. While the focus of this perspective has been predominantly molecule centric, clearly the application of DL outside the pharmaceutical industry attests to the broad array of applications and potential for impact, there is likely overlap in these domains. For example could we learn from social media data using DL methods what might be potential side effects of drugs or even new uses identified by the population. DL could be applied not only to adverse event prediction but also formulation properties, pharmacokinetic simulation, cost effectiveness and even clinical trial simulation and design. How far DL will take us and how quickly it will have an impact on research and the industry remains to be seen, but there is considerable opportunity to develop user accessible tools and apply them now to the accumulating public datasets. DL is already off to a good start in many areas and pharmaceutical researchers would do well to take a closer look and embrace it faster than they have other computational technologies in the past.

Acknowledgments

Dr. Alex M. Clark and Dr. Peter W. Swaan are kindly acknowledged for useful discussions on this topic.

Funding

This work was partially supported by Award Number 9R44TR000942-02 “Biocomputation across distributed private datasets to enhance drug discovery” from the NIH National Center for Advancing Translational Sciences.

Abbreviations

ADME/Tox	Absorption, Distribution, Metabolism, Excretion/Toxicology
AUC	area under the curve
DILI	Drug induced liver injury
hERG	human ether a-go-go related gene
PLGA	poly-lactide-co-glycolide
QSAR	quantitative structure activity relationships

SVM support vector machines**References**

1. Ekins S, Gupta RR, Gifford E, Bunin BA, Waller CL. Chemical space: missing pieces in cheminformatics. *Pharm Res.* 2010; 27(10):2035–2039. [PubMed: 20683645]
2. Rost B, Radivojac P, Bromberg Y. Protein function in precision medicine: deep understanding with machine learning. *FEBS Lett.* 2016
3. Mitchell JB. Machine learning methods in chemoinformatics. *Wiley Interdiscip Rev Comput Mol Sci.* 2014; 4(5):468–481. [PubMed: 25285160]
4. Zhu H, Zhang J, Kim MT, Boison A, Sedykh A, Moran K. Big data in chemical toxicity research: the use of high-throughput screening assays to identify potential toxicants. *Chem Res Toxicol.* 2014; 27(10):1643–1651. [PubMed: 25195622]
5. Clark AM, Ekins S. Open Source Bayesian Models: 2. Mining A “big dataset” to create and validate models with ChEMBL. *J Chem Inf Model.* 2015; 55:1246–1260. [PubMed: 25995041]
6. Ekins S, Clark AM, Swamidass SJ, Litterman N, Williams AJ. Bigger data, collaborative tools and the future of predictive drug discovery. *J Comput Aided Mol Des.* 2014; 28(10):997–1008. [PubMed: 24943138]
7. Ekins S, Freundlich JS, Reynolds RC. Are Bigger Data Sets Better for Machine Learning? Fusing Single-Point and Dual-Event Dose Response Data for Mycobacterium tuberculosis. *J Chem Inf Model.* 2014; 54:2157–2165. [PubMed: 24968215]
8. Ekins S, Ecker GF, Chiba P, Swaan PW. Future directions for drug transporter modelling. *Xenobiotica.* 2007; 37(10):1152–1170. [PubMed: 17968741]
9. Clark AM, Sarker M, Ekins S. New target predictions and visualization tools incorporating open source molecular fingerprints for TB Mobile 2.0. *J Cheminform.* 2014; 6:38. [PubMed: 25302078]
10. Ekins S, Clark AM, Wright SH. Making Transporter Models for Drug-Drug Interaction Prediction Mobile. *Drug Metab Dispos.* 2015; 43:1642–1645. [PubMed: 26199424]
11. Baskin II, Winkler D, Tetko IV. A renaissance of neural networks in drug discovery. *Expert Opin Drug Discov.* 2016; 11:785–795. [PubMed: 27295548]
12. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature.* 2015; 521(7553):436–444. [PubMed: 26017442]
13. Burden F, Winkler D. Bayesian regularization of neural networks. *Methods Mol Biol.* 2008; 458:25–44. [PubMed: 19065804]
14. Gawehn E, Hiss JA, Schneider G. Deep Learning in Drug Discovery. *Mol Inform.* 2016; 35(1):3–14. [PubMed: 27491648]
15. Mamoshina P, Vieira A, Putin E, Zhavoronkov A. Applications of Deep Learning in Biomedicine. *Mol Pharm.* 2016; 13(5):1445–1454. [PubMed: 27007977]
16. Chow, J-F. Things to try after useR! – Part 1: Deep Learning with H2O. Aug 8th. 2016 Available from: <http://www.r-bloggers.com/things-to-try-after-user-part-1-deep-learning-with-h2o/>
17. Anon. TensorFlow. Aug 8th. 2016 Available from: <https://www.tensorflow.org/>
18. Anon. DeepLearning4j. Aug 8th. 2016 Available from: <http://deeplearning4j.org/>
19. Novet, J. Facebook open-sources its cutting-edge deep learning tools. Aug 8th. 2016 Available from: <http://venturebeat.com/2015/01/16/facebook-opens-up-about-more-of-its-cutting-edge-deep-learning-tools/>
20. Chintala, S. FAIR open sources deep-learning modules for Torch. Aug 8th. 2016 Available from: <https://research.facebook.com/blog/fair-open-sources-deep-learning-modules-for-torch/>
21. Linn, A. Microsoft releases CNTK, its open source deep learning toolkit, on GitHub. Aug 8th. 2016 Available from: <http://blogs.microsoft.com/next/2016/01/25/microsoft-releases-cntk-its-open-source-deep-learning-toolkit-on-github/#sm.00013j280xp1sdctrgg21w81es5ov>
22. Angermueller C, Parnamaa T, Parts L, Stegle O. Deep learning for computational biology. *Mol Syst Biol.* 2016; 12(7):878. [PubMed: 27474269]
23. Min S, Lee B, Yoon S. Deep learning in bioinformatics. *Brief Bioinform.* 2016

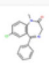


24. Deng X, Gumm J, Karki S, Eickholt J, Cheng J. An Overview of Practical Applications of Protein Disorder Prediction and Drive for Faster, More Accurate Predictions. *Int J Mol Sci.* 2015; 16(7): 15384–15404. [PubMed: 26198229]
25. Akbal-Delibas B, Farhoodi R, Pomplun M, Haspel N. Accurate refinement of docked protein complexes using evolutionary information and deep learning. *J Bioinform Comput Biol.* 2016; 14(3):1642002. [PubMed: 26846813]
26. Zhang S, Zhou J, Hu H, Gong H, Chen L, Cheng C, Zeng J. A deep learning framework for modeling structural features of RNA-binding protein targets. *Nucleic Acids Res.* 2016; 44(4):e32. [PubMed: 26467480]
27. Kandaswamy C, Silva LM, Alexandre LA, Santos JM. High-Content Analysis of Breast Cancer Using Single-Cell Deep Transfer Learning. *J Biomol Screen.* 2016; 21(3):252–259. [PubMed: 26746583]
28. Lusci A, Pollastri G, Baldi P. Deep architectures and deep learning in chemoinformatics: the prediction of aqueous solubility for drug-like molecules. *J Chem Inf Model.* 2013; 53(7):1563–1575. [PubMed: 23795551]
29. Hughes TB, Miller GP, Swamidass SJ. Modeling Epoxidation of Drug-like Molecules with a Deep Machine Learning Network. *ACS Cent Sci.* 2015; 1(4):168–180. [PubMed: 27162970]
30. Aliper A, Plis S, Artemov A, Ulloa A, Mamoshina P, Zhavoronkov A. Deep Learning Applications for Predicting Pharmacological Properties of Drugs and Drug Repurposing Using Transcriptomic Data. *Mol Pharm.* 2016; 13(7):2524–2530. [PubMed: 27200455]
31. Zawbaa HM, Szlek J, Grosan C, Jachowicz R, Mendyk A. Computational Intelligence Modeling of the Macromolecules Release from PLGA Microspheres-Focus on Feature Selection. *PLoS One.* 2016; 11(6):e0157610. [PubMed: 27315205]
32. Ma J, Sheridan RP, Liaw A, Dahl GE, Svetnik V. Deep neural nets as a method for quantitative structure-activity relationships. *J Chem Inf Model.* 2015; 55(2):263–274. [PubMed: 25635324]
33. Xu Y, Dai Z, Chen F, Gao S, Pei J, Lai L. Deep Learning for Drug-Induced Liver Injury. *J Chem Inf Model.* 2015; 55(10):2085–2093. [PubMed: 26437739]
34. Ekins S, Williams AJ, Xu JJ. A Predictive Ligand-Based Bayesian Model for Human Drug Induced Liver Injury. *Drug Metab Dispos.* 2010; 38:2302–2308. [PubMed: 20843939]
35. Ekins S, Litterman NK, Lipinski CA, Bunin BA. Thermodynamic Proxies to Compensate for Biases in Drug Discovery Methods. *Pharm Res.* 2016; 33(1):194–205. [PubMed: 26311555]
36. Perryman AL, Stratton TP, Ekins S, Freundlich JS. Predicting mouse liver microsomal stability with “pruned” machine learning models and public data. *Pharm Res.* 2015; 33:433–449. [PubMed: 26415647]
37. Lienard P, Gavartin J, Boccardi G, Meunier M. Predicting drug substances autoxidation. *Pharm Res.* 2015; 32(1):300–310. [PubMed: 25115828]
38. Fagerberg JH, Karlsson E, Ulander J, Hanisch G, Bergstrom CA. Computational prediction of drug solubility in fasted simulated and aspirated human intestinal fluid. *Pharm Res.* 2015; 32(2):578–589. [PubMed: 25186438]
39. Kingsley LJ, Wilson GL, Essex ME, Lill MA. Combining structure- and ligand-based approaches to improve site of metabolism prediction in CYP2C9 substrates. *Pharm Res.* 2015; 32(3):986–1001. [PubMed: 25208877]
40. Baba H, Takahara J, Mamitsuka H. In Silico Predictions of Human Skin Permeability using Nonlinear Quantitative Structure-Property Relationship Models. *Pharm Res.* 2015; 32(7):2360–2371. [PubMed: 25616540]
41. Baba H, Takahara J, Yamashita F, Hashida M. Modeling and Prediction of Solvent Effect on Human Skin Permeability using Support Vector Regression and Random Forest. *Pharm Res.* 2015; 32(11):3604–3617. [PubMed: 26033768]
42. Wang W, Kim MT, Sedykh A, Zhu H. Developing Enhanced Blood-Brain Barrier Permeability Models: Integrating External Bio-Assay Data in QSAR Modeling. *Pharm Res.* 2015; 32(9):3055–3065. [PubMed: 25862462]
43. Hatanaka T, Yoshida S, Kadhum WR, Todo H, Sugibayashi K. In Silico Estimation of Skin Concentration Following the Dermal Exposure to Chemicals. *Pharm Res.* 2015; 32(12):3965–3974. [PubMed: 26195007]

44. Anon. The Economist. 2016. Special report: The return of the machinery question.
45. Ekins S, Reynolds R, Kim H, Koo M-S, Ekonomidis M, Talaue M, Paget SD, Woolhiser LK, Lenaerts AJ, Bunin BA, Connell N, Freundlich JS. Bayesian Models Leveraging Bioactivity and Cytotoxicity Information for Drug Discovery. *Chem Biol.* 2013; 20:370–378. [PubMed: 23521795]
46. Zhang L, Fourches D, Sedykh A, Zhu H, Golbraikh A, Ekins S, Clark J, Connelly MC, Sigal M, Hodges D, Guiguemde A, Guy RK, Tropsha A. Discovery of novel antimalarial compounds enabled by QSAR-based virtual screening. *J Chem Inf Model.* 2013; 53(2):475–492. [PubMed: 23252936]
47. Mayr A, Klambauer G, Unterthiner T, Hochreiter S. DeepTox: toxicity prediction using deep learning. *Front Environ Sci.* 2016; 3:80.
48. Abdelaziz A, Spahn-Langguth H, Schramm K-W, Tetko IV. Consensus modeling for HTS assays using in silico descriptors calculates the best balanced accuracy in Tox21 challenge. *Front Environ Sci.* 2016; 4:2.
49. Whelan KE, King RD. Intelligent software for laboratory automation. *Trends Biotechnol.* 2004; 22(9):440–445. [PubMed: 15331223]
50. Tetko IV, Bruneau P, Mewes HW, Rohrer DC, Poda GI. Can we estimate the accuracy of ADME-Tox predictions? *Drug Discov Today.* 2006; 11(15-16):700–707. [PubMed: 16846797]
51. Fourches D, Muratov E, Tropsha A. Trust, but verify: on the importance of chemical structure curation in cheminformatics and QSAR modeling research. *J Chem Inf Model.* 2010; 50(7):1189–1204. [PubMed: 20572635]
52. Williams AJ, Ekins S, Tkachenko V. Towards a Gold Standard: Regarding Quality in Public Domain Chemistry Databases and Approaches to Improving the Situation. *Drug Disc Today.* 2012; 17:685–701.
53. Vracko M, Bandelj V, Barbieri P, Benfenati E, Chaudhry Q, Cronin M, Devillers J, Gallegos A, Gini G, Gramatica P, Helma C, Mazzatorta P, Neagu D, Netzeva T, Pavan M, Patlewicz G, Randic M, Tsakovska I, Worth A. Validation of counter propagation neural network models for predictive toxicology according to the OECD principles: a case study. *SAR QSAR Environ Res.* 2006; 17(3): 265–284. [PubMed: 16815767]
54. Ekins S, Wood J. Incentives for Starting Small Companies Focused on Rare and Neglected Diseases. *Pharm Res.* 2016; 33:809–815. [PubMed: 26666772]
55. Ponder EL, Freundlich JS, Sarker M, Ekins S. Computational models for neglected diseases: gaps and opportunities. *Pharm Res.* 2014; 31(2):271–277. [PubMed: 23990313]
56. Murnane, K. What is deep learning and how is it useful? *Forbes.* Available from: <http://www.forbes.com/sites/kevinmurnane/2016/04/01/what-is-deep-learning-and-how-is-it-useful/#715d1eaf10f0>
57. Murnane, K. Thirteen Companies That Use Deep Learning To Produce Actionable Results. *Forbes.* Available from: <http://www.forbes.com/sites/kevinmurnane/2016/04/01/thirteen-companies-that-use-deep-learning-to-produce-actionable-results/#4e710eb07967>
58. Tetko IV, D ML, Williams AJ. The development of models to predict melting and pyrolysis point data associated with several hundred thousand compounds mined from PATENTS. *J Cheminform.* 2016; 8:2. [PubMed: 26807157]
59. Tetko IV, Novotarskyi S, Sushko I, Ivanov V, Petrenko AE, Dieden R, Lebon F, Mathieu B. Development of dimethyl sulfoxide solubility models using 163,000 molecules: using a domain applicability metric to select more reliable predictions. *J Chem Inf Model.* 2013; 53(8):1990–2000. [PubMed: 23855787]
60. Kim S, Thiessen PA, Bolton EE, Chen J, Fu G, Gindulyte A, Han L, He J, He S, Shoemaker BA, Wang J, Yu B, Zhang J, Bryant SH. PubChem Substance and Compound databases. *Nucleic Acids Res.* 2016; 44(D1):D1202–1213. [PubMed: 26400175]
61. Putin E, Mamoshina P, Aliper A, Korzinkin M, Moskalev A, Kolosov A, Ostrovskiy A, Cantor C, Vijg J, Zhavoronkov A. Deep biomarkers of human aging: Application of deep neural networks to biomarker development. *Aging (Albany NY).* 2016; 8(5):1021–1033. [PubMed: 27191382]
62. Di Lena P, Nagata K, Baldi P. Deep architectures for protein contact map prediction. *Bioinformatics.* 2012; 28(19):2449–2457. [PubMed: 22847931]

63. Fakoor, R.; Ladhak, F.; Nazi, A.; Huber, M. Using deep learning to enhance cancer diagnosis and classification; Proceeding of the 30th International conference on machine learning; Atlanta, GA: JMLR: W&CP. 2013;
64. Zeng T, Li R, Mukkamala R, Ye J, Ji S. Deep convolutional neural networks for annotating gene expression patterns in the mouse brain. *BMC Bioinformatics*. 2015; 16:147. [PubMed: 25948335]
65. Chen CL, Mahjoubfar A, Tai LC, Blaby IK, Huang A, Niazi KR, Jalali B. Deep Learning in Label-free Cell Classification. *Sci Rep*. 2016; 6:21471. [PubMed: 26975219]
66. Kraus OZ, Ba JL, Frey BJ. Classifying and segmenting microscopy images with deep multiple instance learning. *Bioinformatics*. 2016; 32(12):i52–i59. [PubMed: 27307644]
67. Park S, Lee SJ, Weiss E, Motai Y. Intra- and Inter-Fractional Variation Prediction of Lung Tumors Using Fuzzy Deep Learning. *IEEE J Transl Eng Health Med*. 2016; 4:4300112. [PubMed: 27170914]
68. Wang C, Liu J, Luo F, Tan Y. Pairwise input neural network for target-ligand interaction prediction. *IEEE Int Conf Bioinf and Biomed*. 2014:67–70.
69. Sushko I, Novotarskyi S, Korner R, Pandey AK, Rupp M, Teetz W, Brandmaier S, Abdelaziz A, Prokopenko VV, Tanchuk VY, Todeschini R, Varnek A, Marcou G, Ertl P, Potemkin V, Grishina M, Gasteiger J, Schwab C, Baskin II, Palyulin VA, Radchenko EV, Welsh WJ, Kholodovych V, Chekmarev D, Cherkasov A, Aires-de-Sousa J, Zhang QY, Bender A, Nigsch F, Patiny L, Williams A, Tkachenko V, Tetko IV. Online chemical modeling environment (OCHEM): web platform for data storage, model development and publishing of chemical information. *J Comput Aided Mol Des*. 2011; 25(6):533–554. [PubMed: 21660515]
70. Walker T, Grulke CM, Pozefsky D, Tropsha A. Chembench: a cheminformatics workbench. *Bioinformatics*. 2010; 26(23):3000–3001. [PubMed: 20889496]
71. Ekins, S.; Hohman, M.; Bunin, BA. Pioneering use of the cloud for development of the collaborative drug discovery (cdd) database. In: Ekins, S.; Hupcey, MAZ.; Williams, AJ., editors. *Collaborative Computational Technologies for Biomedical Research*. Wiley and Sons; Hoboken: 2011.
72. Clark AM, Dole K, Coulon-Spector A, McNutt A, Grass G, Freundlich JS, Reynolds RC, Ekins S. Open source bayesian models: 1. Application to ADME/Tox and drug discovery datasets. *J Chem Inf Model*. 2015; 55:1231–1245. [PubMed: 25994950]
73. Cheng T, Li Q, Wang Y, Bryant SH. Binary classification of aqueous solubility using support vector machines with reduction and recombination feature selection. *J Chem Inf Model*. 2011; 51(2):229–236. [PubMed: 21214224]
74. Low Y, Uehara T, Minowa Y, Yamada H, Ohno Y, Urushidani T, Sedykh A, Muratov E, Kuz'min V, Fourches D, Zhu H, Rusyn I, Tropsha A. Predicting drug-induced hepatotoxicity using QSAR and toxicogenomics approaches. *Chem Res Toxicol*. 2011; 24(8):1251–1262. [PubMed: 21699217]
75. Ekins S. Progress in computational toxicology. *J Pharmacol Toxicol Methods*. 2014; 69(2):115–140. [PubMed: 24361690]
76. Wang S, Li Y, Wang J, Chen L, Zhang L, Yu H, Hou T. ADMET evaluation in drug discovery. 12. Development of binary classification models for prediction of hERG potassium channel blockage. *Mol Pharm*. 2012; 9(4):996–1010. [PubMed: 22380484]
77. Kortagere S, Chekmarev DS, Welsh WJ, Ekins S. New predictive models for blood brain barrier permeability of drug-like molecules. *Pharm Res*. 2008; 25:1836–1845. [PubMed: 18415049]
78. Leong MK. A novel approach using pharmacophore ensemble/support vector machine (PhE/SVM) for prediction of hERG liability. *Chem Res Toxicol*. 2007; 20(2):217–226. [PubMed: 17261034]
79. Hou T, Wang J, Li Y. ADME evaluation in drug discovery. 8. The prediction of human intestinal absorption by a support vector machine. *J Chem Inf Model*. 2007; 47(6):2408–2415. [PubMed: 17929911]
80. Clark AM, Dole K, Ekins S. Open Source Bayesian Models: 3. Composite Models for prediction of binned responses. *J Chem Inf Model*. 2016; 56:275–285. *J Chem Inf Model*. [PubMed: 26750305]
81. Kim S, Jin D, Lee H. Predicting drug-target interactions using drug-drug interactions. *PLoS One*. 2013; 8(11):e80129. [PubMed: 24278248]

82. Unterthiner, T.; Mayr, A.; Klambauer, G.; Hochreiter, S. Toxicity prediction using deep learning. Available from: <https://arxiv.org/pdf/1503.01445.pdf>
83. Zheng X, Ekins S, Raufman JP, Polli JE. Computational models for drug inhibition of the human apical sodium-dependent bile acid transporter. *Mol Pharm*. 2009; 6(5):1591–1603. [PubMed: 19673539]
84. Diao L, Ekins S, Polli JE. Quantitative Structure Activity Relationship for Inhibition of Human Organic Cation/Carnitine Transporter. *Mol Pharm*. 2010; 7:2120–2130. [PubMed: 20831193]
85. Dong Z, Ekins S, Polli JE. Structure-activity relationship for FDA approved drugs as inhibitors of the human sodium taurocholate cotransporting polypeptide (NTCP). *Mol Pharm*. 2013; 10(3): 1008–1019. [PubMed: 23339484]
86. You H, Lee K, Lee S, Hwang SB, Kim KY, Cho KH, No KT. Computational classification models for predicting the interaction of compounds with hepatic organic ion importers. *Drug Metab Pharmacokinet*. 2015; 30(5):347–351. [PubMed: 26293543]
87. de Cerqueira Lima P, Golbraikh A, Oloff S, Xiao Y, Tropsha A. Combinatorial QSAR modeling of P-glycoprotein substrates. *J Chem Inf Model*. 2006; 46(3):1245–1254. [PubMed: 16711744]
88. Xue Y, Yap CW, Sun LZ, Cao ZW, Wang JF, Chen YZ. Prediction of P-glycoprotein substrates by a support vector machine approach. *J Chem Inf Comput Sci*. 2004; 44(4):1497–1505. [PubMed: 15272858]
89. Xu C, Cheng F, Chen L, Du Z, Li W, Liu G, Lee PW, Tang Y. In silico prediction of chemical Ames mutagenicity. *J Chem Inf Model*. 2012; 52(11):2840–2847. [PubMed: 23030379]
90. Hansen K, Mika S, Schroeter T, Sutter A, ter Laak A, Steger-Hartmann T, Heinrich N, Muller KR. Benchmark data set for in silico prediction of Ames mutagenicity. *J Chem Inf Model*. 2009; 49(9): 2077–2081. [PubMed: 19702240]
91. Moss GP, Shah AJ, Adams RG, Davey N, Wilkinson SC, Pugh WJ, Sun Y. The application of discriminant analysis and Machine Learning methods as tools to identify and classify compounds with potential as transdermal enhancers. *Eur J Pharm Sci*. 2012; 45(1-2):116–127. [PubMed: 22101136]
92. Vock DM, Wolfson J, Bandyopadhyay S, Adomavicius G, Johnson PE, Vazquez-Benitez G, O'Connor PJ. Adapting machine learning techniques to censored time-to-event health record data: A general-purpose approach using inverse probability of censoring weighting. *J Biomed Inform*. 2016; 61:119–131. [PubMed: 26992568]
93. Chia CC, Rubinfeld I, Scirica BM, McMillan S, Gurm HS, Syed Z. Looking beyond historical patient outcomes to improve clinical models. *Sci Transl Med*. 2012; 4(131):131ra149.
94. Rochefort CM, Verma AD, Egual T, Lee TC, Buckeridge DL. A novel method of adverse event detection can accurately identify venous thromboembolisms (VTEs) from narrative electronic health record data. *J Am Med Inform Assoc*. 2015; 22(1):155–165. [PubMed: 25332356]
95. Degardin K, Guillemain A, Guerreiro NV, Roggo Y. Near infrared spectroscopy for counterfeit detection using a large database of pharmaceutical tablets. *J Pharm Biomed Anal*. 2016; 128:89–97. [PubMed: 27236101]
96. Khamis MA, Gomaa W, Ahmed WF. Machine learning in computational docking. *Artif Intell Med*. 2015; 63(3):135–152. [PubMed: 25724101]

SAR data
+ Molecular descriptors

Ingredient	Molecular Weight	Standard Type	Standard Value	Standard Units	Assay Type
 CHEM12	284.34	PPH	87.82	%	A
 CHEM18	252.27	PPH	82.88	%	A
 CHEM388	454.6	PPH	80.88	%	A

Calculated Compound Parent Properties						
Mol. Weight	Mol. Weight Monoisotopic	ALogP	#Rotatable Bonds	Polar Surface Area	Molecular Species	
252.3	252.0899	2.79	2	61.69	ZWITTERION	

HBA	HBD	#Ro5 Violations	HBA (Lipinski)	HBD (Lipinski)	#Ro5 Violations (Lipinski)
4	2	0	4	2	0

ACD Acidic pKa	ACD Basic pKa	ACD LogP	ACD LogD pH7.4	Aromatic Rings	Heavy Atoms	QED Weighted
1.84	10.79	1.74	-.71	2	19	0.86

Input Layer

Hidden Layers

Output layer

Correlation

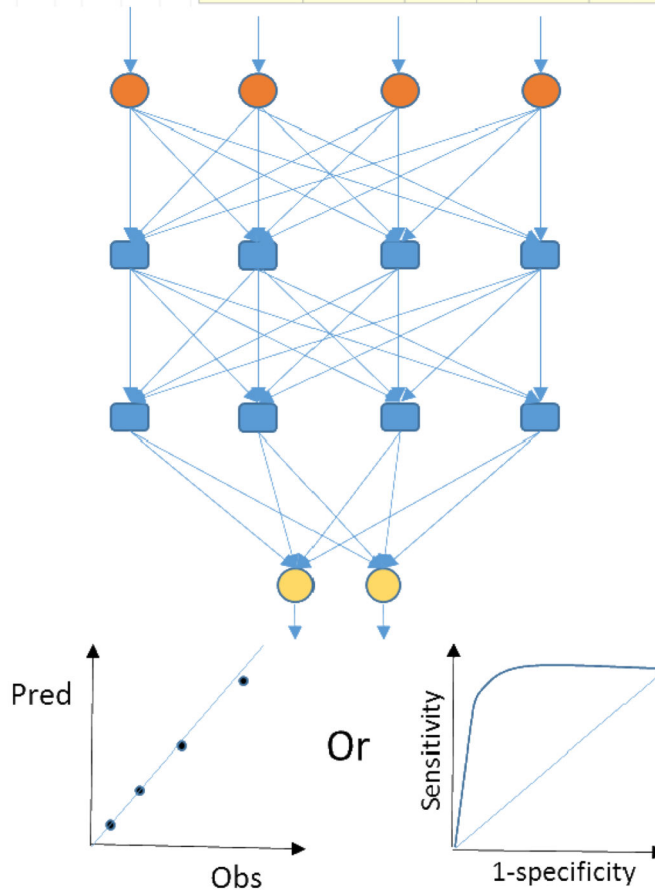


Figure 1.

A schematic of a deep learning neural network applied to cheminformatics and a single property with output as a quantitative or qualitative prediction.

Public and private
databases and molecular
descriptors

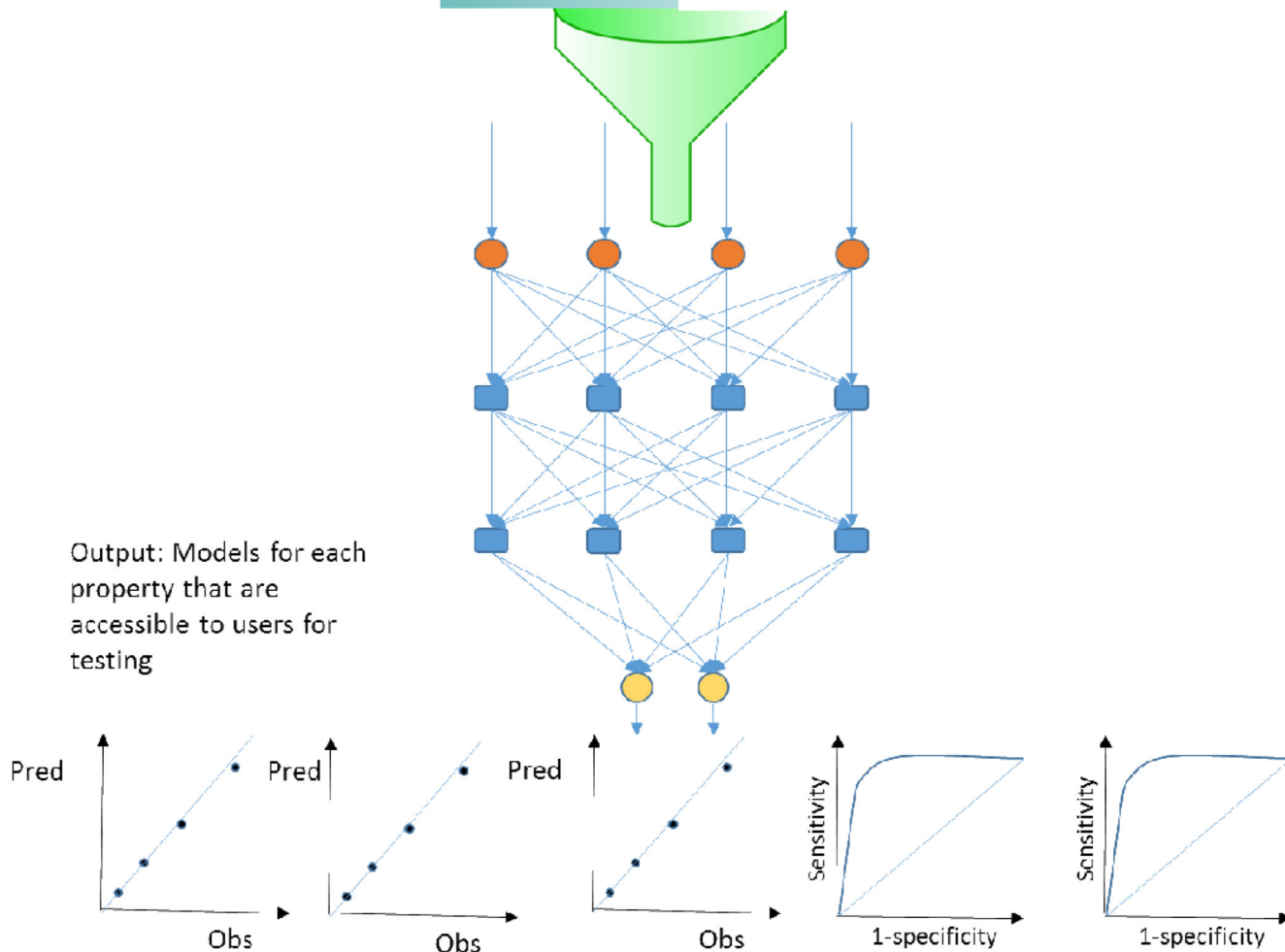


Figure 2.

Using public and private data for generating deep learning models for application across vast numbers of endpoints whether quantitative or qualitative (classification models).

Table 1

Examples of deep learning studies. Deep learning (DL), Support Vector Machine (SVM) and artificial neural network (ANN)

Example end point modelled	Dataset size	Model statistics	Summary	References
Solubility	1144, 1026, 74, 125	0.92 RMSE 0.58, 0.91 RMSE 0.60, 0.81 RMSE 0.72, 0.67 RMSE 0.90 All 10 fold cross validation	Did not compare to other machine learning methods themselves. Addition of log P in some cases improved models.	(28)
Drug Induced Liver Injury	190, 475, 1065	80.5, 88.4, 70.1, Accuracies from internal cross validation	DL Models also assessed with external test sets. DL outperformed ANN with different PaDEL and Mold2 descriptors used.	(33)
ADME and target activity	Multiple models from 2092 (microsomes) to 318,795 (hERG)	Compared DL to random forest (RF) models on external datasets (11 of 15 DL models outperform RF models from Kaggle test and 13 out of 15 additional models)	15 Kaggle datasets available consisting of activities and descriptors.	(32)
Biomarkers	62,419 records and 46 blood markers	$R^2 = 0.80$ and 82% prediction accuracy in predicting chronological age within a 10 year window.	Ensemble model performs better $R^2 = 0.82$ accuracy 83.5%	(61)
Protein contact maps	ASTRAL database release 1.75	CASP 8 and 9 data is used for test sets. CMAPpro more accurate on the CASP 8 dataset than other methods. Improvement 10% or higher.	Predicts contact maps with accuracy of 30%	(62)
Cancer diagnosis	13 published cancer gene expression datasets	10 fold cross validation average classification accuracy 46.33 – 100%	Principal Component Analysis was used to preprocess data	(63)
Gene expression patterns	2000 genes from Allen Developing Mouse Brain Atlas	Average AUC 0.894 versus 0.82 for bag of words for annotating gene expression pattern. Gene ontology functional annotation average AUC =	Used deep convolutional neural network. Relative performance of different methods differs across different developing stages.	(64)

Example end point modelled	Dataset size	Model statistics	Summary	References
		0.59 vs 0.57 for bag of words		
Protein disorder prediction	1111 proteins	AUC 86.8 on training set AUC 80.9 on CASP10 dataset	More sophisticated methods such as DISOPRED3 and DNDISORDER have test AUC 87.2 and 82.3 respectively	(24)
RNA binding protein features	24 datasets	10 fold cross validation used AUROC values 0.71 – 0.99	Outperformed the GraphProt method when using base sequence, secondary and tertiary structural profiles.	(26)
High content analysis of breast cancer	148,649 rows	Leave one out cross validation Accuracy 87%	Linear SVM Accuracy 20.95, SVM using radial basis function Accuracy 21.04 Uses a public dataset for training.	(27)
Epoxidation	702	Site of epoxidation AUC 94.9% and separation of epoxidation vs non epoxidation molecules with AUC 79.3% after leave one out cross validation	The DL was compared with a logistic regression model which gave Site of epoxidation AUC 93.7% and epoxidation vs non epoxidation molecules AUC 78.9%	(29)
Tox21	11,764 training set, 296 Leaderboard set, 647 test set	DL average AUC = 0.837, SVM average AUC = 0.832, RF average AUC = 0.803	DL with multitask learning outperformed single task learning on 10 of 12 assays and DL won 9 of 15 challenges with nuclear receptor and stress response panels	(47)
Tox21	11,764 training set, 296 Leaderboard set, 647 test set	Consensus model Balanced accuracies from 0.599-0.903	Analyzed 12 targets using associative neural networks. Built consensus models. Training and leaderboard sets were combined. Reported to have the best balanced accuracy	(48)
Refinement of docked proteins	35,000 samples of 35 unbound dimer proteins	Tested on 25 test cases across 5 proteins - RMSD 1.40Å	No information on dataset availability	

Example end point modelled	Dataset size	Model statistics	Summary	References
Drug repurposing	977 Landmark genes 271 signaling pathways	10 fold cross validation deep learning outperforms SVM for both datasets. DL F1 score of 0.70 and SVM 0.53 for pathway 3 class problem.	Deep learning models trained just on gene data did not perform well at classifying 12 groups of drugs. Pathways performed better.	(30)
Label free cell classification	Not defined	5 fold Cross validation shows 16 multivariate features can provide a balanced accuracy of 96.4%	Classification of blood cells (OT-II) and cancer cells (SW-480). Deep learning outperforms logistic regression, SVM and Bayesian but all accuracies are greater than 85%	(65)
Classifying microscopy images	103 treatments – 25 images per class for training	Test set accuracy 0.96	Used an available MFC-7 breast cancer imageset BBBC021v1	(66)
Lung tumors	130 patients	Improved mean square error by 29.98% and prediction overshoot by 70.93%	Computation time of 1.54ms might achieve real-time estimation of intra-fractional variation and better tracking for radiotherapy.	(67)
Target-ligand interaction prediction	sc-PDB, 836 targets and 2710 ligands	5 fold cross validation AUC 0.959. Outperforms other published methods such as BLM-NII (AUC 0.858) and CS-PD (AUC 0.799)	Used a pairwise input neural network. % fold accuracy = 0.887	(68)

Table 2

Representative examples of machine learning models applied to pharmaceutically relevant end points to indicate areas used for machine learning that could be useful datasets for potential future use of deep learning. Other publically accessible datasets are available in OCHEM (69), Chembench (70), CDD (71) etc.

Example end point modelled	Naïve Bayesian	Support Vector Machine	Deep Learning
Solubility	(72)	(73)	(28)
Drug Induced Liver Injury	(34)	(74)	(33)
hERG	(75, 76)	(77, 78)	(32)
ADME	(36, 72)	(79)	(29, 32)
Blood Brain Barrier penetration	(75)	(77)	
Biological Targets	(5, 80)	(81)	(32, 82)
Skin Permeability		(40, 41)	
Transporters	(10, 83-85)	(86-88)	(32)
Mutagenicity	(75, 89)	(89, 90)	
Formulation		(91)	
Adverse event prediction	(92)	(93, 94)	
Counterfeit drug detection		(95)	
Docking		(96)	
Small molecule pKa		(3)	