



Published in final edited form as:

Chem Soc Rev. 2020 June 07; 49(11): 3525–3564. doi:10.1039/d0cs00098a.

QSAR without borders

Eugene N. Muratov^{a,b}, Jürgen Bajorath^c, Robert P. Sheridan^d, Igor Tetko^e, Dmitry Filimonov^f, Vladimir Poroikov^f, Tudor I. Oprea^g, Igor I. Baskin^{h,i}, Alexandre Varnek^h, Adrian Roitberg^j, Olexandr Isayev^a, Stefano Curtarolo^k, Denis Fourches^l, Yoram Cohen^m, Alan Aspuru-Guzikⁿ, David A. Winkler^o, Dimitris Agrafiotis^p, Artem Cherkasov^{*,q}, Alexander Tropsha^{*,a}

^aUNC Eshelman School of Pharmacy, University of North Carolina, Chapel Hill, NC, USA.

^bDepartment of Pharmaceutical Sciences, Federal University of Paraiba, Joao Pessoa, PB, Brazil.

^cDepartment of Life Science Informatics, University of Bonn, Bonn, Germany.

^dMerck & Co. Inc, Kenilworth, NJ, USA.

^eInstitute of Structural Biology, Helmholtz Zentrum München - Deutsches Forschungszentrum für Gesundheit und Umwelt (GmbH) and BIGCHEM GmbH, Neuherberg, Germany.

^fInstitute of Biomedical Chemistry, Moscow, Russia.

^gDepartment of Internal Medicine and UNM Comprehensive Cancer Center, University of New Mexico, Albuquerque, NM, USA; Department of Rheumatology, Gothenburg University, Sweden; Novo Nordisk Foundation Center for Protein Research, University of Copenhagen, Denmark

^hDepartment of Chemistry, University of Strasbourg, Strasbourg, France.

ⁱFaculty of Physics, M.V. Lomonosov Moscow State University, Moscow, Russia

^jDepartment of Chemistry, University of Florida, Gainesville, FL, USA.

^kMaterials Science, Center for Autonomous Materials Design, Duke University, Durham, NC, USA.

^lDepartment of Chemistry, North Carolina State University, Raleigh, NC, USA.

^mInstitute of The Environment and Sustainability, University of California, Los Angeles, CA, USA

ⁿDepartment of Chemistry, University of Toronto, Toronto, ON, Canada.

^oMonash Institute of Pharmaceutical Sciences, Monash University, Melbourne, VIC, Australia; La Trobe Institute for Molecular Science, La Trobe University, Bundoora, Australia; CSIRO Manufacturing, Clayton, Australia; School of Pharmacy, University of Nottingham, Nottingham, UK.

^pNovartis Institutes for BioMedical Research (NIBR), Cambridge, MA, USA.

* corresponding authors.

Conflicts of interest

There are no conflicts to declare.

^qVancouver Prostate Centre, University of British Columbia, Vancouver, BC, Canada

Abstract

Prediction of chemical bioactivity and physical properties has been one of the most important applications of statistical and more recently, machine learning and artificial intelligence methods in chemical sciences. This field of research, broadly known as Quantitative Structure-Activity Relationships (QSAR) modeling, has developed many important algorithms and has found a broad range of applications in physical organic and medicinal chemistry in the past 55+ years. This Perspective summarizes recent technological advances in QSAR modeling. Importantly, it also highlights the applicability of algorithms, modeling methods, and validation practices developed in QSAR to a wide range of research areas beyond traditional QSAR fields. These fields include nanotechnology, materials science, biomaterials, clinical informatics, and others. As modern research methods generate rapidly increasing amounts of data, knowledge of robust data-driven modelling methods is becoming essential for scientists in many disciplines both within and outside of chemical research. We hope that this contribution will serve to address this challenge.

Introduction

Quantitative Structure-Activity Relationship (QSAR) modeling is a well-established computational approach to chemical data analysis. QSAR models are developed by establishing empirical, linear or non-linear relationships between values of *chemical descriptors* computed from molecular structure and experimentally measured properties or bioactivities of those molecules, followed by application of these models to predict or design novel chemicals with desired properties.

Historically, QSAR modeling has largely been applied to computer-aided drug discovery. Many papers, reviews, and book chapters describing the methods and applications of QSAR modeling have appeared in the scientific literature since the seminal publication by Hansch et al. in 1962¹ that effectively pioneered the field. More than five years ago, some of the contributors to this paper coauthored a comprehensive review of QSAR modeling,² where we discussed the evolution of methods and best practices of QSAR. The field has grown and evolved substantially subsequently. The Web of Science core collection lists more than 5600 papers on QSAR published within last five years, a substantial fraction of the ~20,000 papers that have been published on this subject since 1962. Many publications have advanced the traditional areas of QSAR modeling such as prediction of biological activities and ADME/Tox properties, building on successful use of QSAR modeling in chemical, agrochemical, pharmaceutical³, and cosmetic industries.⁴ However, new and interesting directions and application areas have also emerged, such as process chemistry^{5,6} and (retro)synthetic route prediction and optimization.⁷ Thus, models have become an integral component of the drug discovery process, providing substantial guidance in planning experiments.^{3,8}

Clearly, QSAR modeling is an established and useful computational chemistry approach. However, many practitioners still consider it limited to modeling and prediction of chemical bioactivities and/or properties. One aim of this Perspective is to outline the opportunities

presented by recent and emerging developments in artificial intelligence (AI), machine learning (ML) and other approaches to modeling Big Data *within* traditional QSAR. However, our prime objective is to emphasize the impact that QSAR methods and approaches have, or will shortly have, on many modern data-driven areas of molecular research *beyond* traditional QSAR areas.

In cheminformatics molecules are represented by mathematical descriptors that encode molecular structures and properties. Multivariate statistical methods or machine learning are employed to establish relationships between descriptors and a target property, such as molecular bioactivity. It is easy to see that analogous representations can be generated for many types of data where *objects* are represented by their *features*, and the general objective is to predict object *properties* (endpoints) from these features. For instance, in clinical data, the objects would be patients, the features would be clinical or pharmacological biomarkers characteristic of the patients, and the target property would be the any health outcomes such as the rate of patient survival.

Regardless of the nature of the data, the same machine learning (ML) approaches can be used universally to analyze and process data in any domain. Furthermore, despite differences in the information content and meaning of the data, different research fields share similar data handling routines. These often replicate the workflows and protocols already created, evaluated, and used in QSAR. Indeed, the general data cycle associated with QSAR projects (Figure 1) can be easily adopted for similar data-analytical investigations in other fields. To further illustrate this point, Table 1 provides a collection of recent references describing studies in diverse research areas that cite some or many concepts from QSAR. Examples include fields as diverse as climatology,⁹ urban engineering,¹⁰ student admissions,¹¹ remote sensing¹² and clinical informatics (discussed in one of the sections of this contribution). Importantly, QSAR modeling was one of the first research fields that highlighted the importance of data curation,¹³ rigorous validation of developed models,¹⁴ and data reproducibility,¹⁵ that has recently become a significant concern to the scientific community.¹⁶

Here we integrate contributions from some of the leading experts in QSAR modeling that illustrate the breadth and generality of modern data processing and modeling practices in the field.¹⁷ The contributors have worked both on methodology and applications of QSAR modeling for most of their professional life. Some of the coauthors have pivoted their research into other areas where QSAR-like approaches have not been used before, illustrating the main theme of this paper by their own careers. We engaged other scientists who work in areas where data modeling was not common but who have started using QSAR-like methods in their research. We are confident that many fields that employ statistical modeling approaches will benefit significantly from the experience accumulated within the QSAR community in the last 55 years.

We start this contribution by discussing fundamental concepts of QSAR, such as chemical similarity. We describe the impact of recent advances, such as Deep Learning (DL), on traditional areas of QSAR modeling, such as drug discovery and development and chemical safety prediction. We then reflect how the complexity of algorithms and the size, diversity,

and complexity of chemical bioactivity data have grown. We also illustrate how modern computational methods are capable of modeling multiple bioactivity endpoints simultaneously, addressing the issue of multi-objective optimization. We then extend traditional boundaries of QSAR by summarizing recent, exciting developments in organic synthesis planning and retrosynthetic pathway prediction, advances in robotic chemistry, and applications of machine learning to quantum chemistry. Finally, to further illustrate the breadth of applicability of modern QSAR approaches, we discuss their use in materials and nanomaterials science, regenerative medicine, and health care. Throughout the discussion, we identify methodological similarities between drug discovery approaches and those employed in other areas. We further propose that experience and best practice of data curation, model development, and validation accumulated by the QSAR community provides valuable guidance for many areas where statistical and machine learning data modeling is applied.

This broad, platform applicability of QSAR algorithms and protocols across all data-rich areas of modern science underpins the appeal of QSAR as a robust, predictive data analysis and modelling tool. We advise contemporary chemists to become familiar with the major computational approaches discussed in this contribution. To this end, borrowing from a recent “In the Pipeline” blog by Derek Lowe,¹⁸ “*it is not that machines are going to replace chemists. It’s that the chemists who use machines will replace those that don’t!*” We hope that this paper will stimulate experimental scientists to consider deeper integration of computational methods and models into their research projects, to consider how the data they generate will be modelled when planning experiments and will serve as useful reference for computational chemists as well.

Chemical similarity

Classical QSAR is defined by linear (regression) models derived from a set of small molecules sharing the same (target-specific) biological activity. A QSAR model predicts changes in potency as a function of structural modifications.^{1,19} The evolution of QSAR modeling from linear to more complex *machine learning* models addressing non-linear relationships between chemical structure and bioactivity was discussed in a paper co-written by one of the founders of classical QSAR, Prof. Toshio Fujita in 2016.¹⁹ Chemical bioactivity data employed in model development are generally derived from investigations of analog series from medicinal chemistry. These sets of compounds usually share a common core structure (scaffold) and carry different substituents (R-groups) at one or more sites. Descriptor-based linear regression models then predict potency of newly designed analogs to further extend such *congeneric* series, a fundamental task of classical QSAR. This prediction scheme is provides a useful guide to compound design and synthesis, making QSAR one of the most popular predictive approaches in medicinal chemistry since its seminal development.¹

QSAR modeling is based upon the premise that structurally similar compounds exhibit similar biological effects, often referred to as the *similarity-property principle* (SPP) The SPP postulates a causal link between molecular similarity and biological activity, which implies that gradual changes in chemical structure are accompanied by gradual changes in

activity. This expectation provides the initial rationale for the generation of linear QSAR models.¹

Chemical similarity is often evaluated in relation to bioactivity. Multi-dimensional structure-activity relationship (SAR) landscapes derived from models, describe similarity relationships between active molecules and their biological potency differences. These can be used to understand the effects of various structural features on biology, especially SAR continuities versus discontinuities in compound responses.²⁰ SAR continuity is directly associated with the SPP, implicating a smooth continuous relationship between conservative structural modifications of active compounds and accompanying moderate potency alterations. In contrast, SAR discontinuities²¹ occur when small structural modifications lead to very large biological potency changes, not consistent with the SPP and falling outside the applicability domain of linear QSAR models. Figure 2 shows small sets of active compounds that are characterized by SAR continuity and discontinuity, respectively. “Activity cliffs” are formed by analogs displaying the largest potency differences in a compound series for the smallest change in structure.²² The existence of activity cliffs in compound data sets is a major factor limiting QSAR predictions, often much greater than intrinsic limitations of modeling.²² Strikingly similar observations have also been made in bioinformatics where some pairs of proteins with high sequence similarity possess very different structures and functions.²³ This analogy is one of many that methodologically bridge between QSAR and other fields that rely on data analytics. It should be noted that activity cliffs may be sensitive to both the choice of descriptors and the degree of the experimental variability. Importantly, SAR discontinuity limits QSAR modeling regardless of molecular representations and descriptors that are used when the corresponding compounds are close structural analogs. Activity landscapes of compound data sets might be “flattened” by using large numbers of features as molecular representations such that compounds become increasingly dissimilar (i.e., their distances in feature space increase). However, introducing artificial dissimilarity results in a loss of SAR information (and often leads to overfitting of regression models).

In QSAR modeling the presence of SAR continuities and discontinuities in sets of active compounds is not mutually exclusive. Rather, continuous and discontinuous SARs coexist in many data sets²¹ resulting in the presence of adjacent gently sloped and rugged regions in activity landscapes (Figure 2). Focusing potency predictions around local regions of SAR continuity can often lead to QSAR models with high predictive power. To this end, numerical SAR analysis methods can be used to identify compound subsets having desirable SAR characteristics.²⁴

Going beyond the traditional QSAR paradigm means departing from the SPP. Modeling compounds with increasingly diverse structures with few or no common scaffolds means that structural differences between active compounds are not gradual, such as those that arise from “scaffold hopping”.²⁵ This leads to structurally diverse active compounds that require non-linear approaches to modeling SARs satisfactorily, making bioactivity predictions more difficult. Non-linear SAR models require analysis of relationships between structure of both close and remote structural analogs and respective changes in their potency. This is beyond

the capacity of classical linear regression QSAR methods and generally requires the use of machine learning (ML) as discussed in the next section.²⁶

To summarize, the choice of molecular representations (descriptors) and assessment of molecular similarity play a critical role in QSAR. It should be emphasized that comparison of object representations, their similarity metrics and the interplay between object relationships and associated (latent) properties is of general relevance for data modeling irrespective of research areas. In fact, the *similiar similibus curantur* (“likes are cured by likes”) principle formulated by Paracelsus²⁷ (the “father of toxicology”) could be seen as one of the most common ways of rational thinking (reflected in the SPP principle as applied in cheminformatics) and reasoning approaches in nearly any area of science. As highlighted throughout this contribution, this principle is one of key drivers of the general applicability of approaches and tools employed in cheminformatics.

Modern trends in QSAR modeling

Chemical similarity may help with qualitative assessment of compound bioactivity but its quantitative evaluation requires the use of statistical tools that can model the relationship between chemical structure and bioactivity.¹ Currently, there is much talk about the use of artificial intelligence (AI) in chemistry. Here we distinguish between AI and machine learning in the following way. AI is the superset of tasks that demonstrate characteristics of human intelligence, while ML is a subset of AI which accesses data, analyses trends and generates intelligent, actionable insights. Many people use the term AI in the same context as ML in many data-rich disciplines, ranging from health care to astronomy. In this regard one can say that AI has been used in chemistry since the 1960’s under the name QSAR. In general, ML represents a set of techniques for predicting a property Y based on known examples, where each example i has property Y(i) and a set of k features X(i,j), j=1 to k. In this section we show how QSAR modeling can be applied much more broadly than has been the case previously. Theoretical organic chemistry, a highly specialized field, gave rise to the QSAR paradigm. The experience and trends in modern QSAR we summarize in this section is illustrative, and perhaps, instructional, for any data-rich area of research.

Machine learning suffers from the same philosophical limitations that any type of inductive learning does: distinguishing correlation from causation and knowing when we have enough training examples to generate a mode that makes accurate predictions for new cases, etc. In QSAR, the dependent variable Y is usually some biological or physical property, and the independent variable features X (descriptors in chemistry) are derivable from chemical structures. In QSAR, historically the objects are drug-sized molecules, but that is not always the case. Objects can be atoms, protein sequences, pairs of proteins, etc., so long as relevant descriptors can be generated.

Chemical descriptors for drug-sized molecules fall into two main categories: substructures, which note the presence and/or frequency of certain groups, and computable properties that are representative of the entire molecule. In QSAR, the function that maps Y from X is called a model. Obviously, the same general construct is used in statistical modeling in any field, except the nature of descriptors depends on the type of the objects.

This section concentrates on trends in QSAR in the pharmaceutical industry because, arguably, that is where the opportunities and challenges for innovation and potential impact on society are greatest.²⁸ Most pharmaceutical companies are likely to develop QSAR models for on-target (e.g., binding of ligands to targets) and off-target (secondary pharmacology) activities, as well as ADMET (absorption, distribution, metabolism, excretion, and toxicity) properties, which are discussed in the next section. Companies also develop their own best practices for building and using QSAR models. Models are used so that predictions can substitute for experiment under some circumstances. However, the current state of the art in QSAR modeling often precludes chemists from relying fully on individual quantitative predictions, rather on predictions of trends accurate enough to prioritize sets of compounds for synthesis and experimental evaluation.

Researchers are always seeking ways to improve their science, and the field of QSAR is no exception. There are many recent trends but here we describe the most important ones that in our opinion, can be generalized to many other research fields:

1. **Data.** Data driven modeling methods are clearly highly dependent on data size, quality, and diversity. The size and diversity of datasets have dramatically increased in recent years due to technological advances in robotics and miniaturization (similar trends of course are observed in nearly any area of research and technology development). We can now generate very large volumes of data for a specific project, typically for 10^4 - 10^6 diverse molecules. Data generation is resource intensive, and data are always containing experimental error. Outside of the pharmaceutical industry, the availability of large volumes of published, or otherwise public domain data in databases like ChEMBL²⁹ and PubChem³⁰ has transformed the field.
2. **Validation methods.** A common method of validating a QSAR model is by use of an external test set. Part of the data is held aside and the remainder used to train the model. The model is used to predict the test set endpoints and a metric for the accuracy of prediction is then calculated. A better way to simulate the natural evolution of a typical drug discovery project is to use a time-split test set,³¹ i.e., assigning compounds tested in later phases of the project to the test set. It can be demonstrated that time-split gives a good estimate of the R^2 for true prospective prediction relative to random test set selection (a standard method that can overestimate prediction accuracy) and leave-class-out validation (which is too pessimistic).³¹ Users of the ChEMBL database sometimes use the date of publication as a surrogate time-split threshold. Validation of QSAR models for properties of chemical mixtures is more complicated. In that regard, the points out³² approach is not different from traditional QSAR, but should be used only for predicting the same mixtures with new composition. The compounds out³² approach is suitable for predicting new mixtures of compounds from the modeling set; the mixtures out³³ approach is for mixtures of one compound from the modeling set and one new compound; and the everything out³⁴ approach (the most rigorous) is for mixtures of completely new compounds.

- 3. Multitask modeling.** In classical QSAR only one predicted activity is modelled at a time. However, in drug development, multiple activities, both on- and off-target, are needed for prioritizing compounds. The set of techniques for prioritizing compounds based on more than one predicted activity simultaneously is called *multi-parameter optimization*,³⁵ or multi-task modeling. In general, this objective can be achieved by an ensemble of single task models, or by a single model that can predict more than one activity simultaneously using either non-neural net or neural net-based techniques, including deep learning that has become popular in recent years. The multiple activities could involve related targets in one species, the same target in different species, the same target under different experimental conditions, or be completely unrelated. Multitask modeling is expected to be useful when data are sparse, i.e. not all molecules are tested on all targets, and the hope is that information will “leak” or “read across” different targets and reinforce structure-activity trends. Several methods have been proposed for multitask QSAR modeling including Perturbation Theory + Machine Learning (PTML),³⁶ inductive learning and multi-objective optimization³⁷ as applied in proteochemometrics modeling.³⁸ The most common way of handling multitask modeling currently is with deep neural nets, especially convolutional neural nets. This will be discussed in more detail in the section on ML methods. Multi-task optimization represents an active area of development in QSAR modeling. However, it is still unclear whether these techniques provide a significant improvement in external predictive accuracy compared to an ensemble of single task models developed for the same end points. For example, an ensemble of models developed with XGBoost (gradient boosting decision trees) method exhibited the best performance in a recent 2019 IDG-DREAM Drug-Kinase Binding Prediction Challenge.³⁹ As many compounds do have multiple biological activities, there is an obvious need to continue both methodological and application studies on multitask modeling in QSAR and other areas of statistical data analysis.
- 4. Applicability Domain (AD).** An applicability domain⁴⁰ defines the space of molecular features on which the model has been trained and to which it should be applied; The AD provides a means for estimating the reliability of property predictions for new molecules from a QSAR model. It allows flagging of less reliable predictions and helps identify additional molecules that might be required to expand the model AD into more productive chemical spaces. Interestingly, AD is one area where QSAR is ahead of the general field of ML, although there is not yet a consensus on the best approach to this issue.⁴⁰
- 5. Modelability.** Whether a statistically significant model can be built from a given dataset depends on a number of issues. If the size of the experimental error in the measured dependent variable approaches the magnitude of the variation across multiple molecules in the dataset, it becomes increasingly hard to generate meaningful models. The signal to noise ratio in the data set is too low. Assuming this is not an issue, and considering activity and descriptors together, the relatively new concept of *modelability*⁴¹ proposes that predictivity of QSAR

models is then limited by *activity cliffs*. As discussed above, activity cliffs exist when very similar compounds have very different activities, making the target property of compounds near the activity cliffs hard to predict.²² This difficulty is not easily overcome by changing either the QSAR method or the descriptors used. One exception is that using stereochemically-aware descriptors can reduce activity cliffs where different stereoisomers exhibit very different activities. Metrics that measure the prevalence of activity cliffs in a dataset are good predictors of the modelability of that dataset.⁴¹ Clearly, these metrics cannot distinguish activity cliffs that are intrinsic to the SAR response surface from those that are artifacts due to large experimental uncertainties in the measured activities.

- 6. Interpretability.** Early classical QSAR methods were relatively simple and tended to deal with molecules that were close analogs. Comparative Molecular Field Analysis (CoMFA)⁴² was extremely successful because of its visual appeal – it was clear where and how to modify a molecule to increase its activity. Later, projection of atom/fragment model contributions onto exemplar molecules has been suggested.⁴³ However, as modeling methods have become more sophisticated, descriptors more arcane, and datasets more diverse, the accuracy and breadth of predictions have increased at the expense of interpretability (understanding the molecular basis for good or bad activity of molecules that guides design of improved examples). Methods that “see” into the black box of QSAR models independent of the descriptors and QSAR methods used are discussed in a recent review.⁴⁴ An important process in QSAR modeling is selecting the most relevant subset of descriptors for a much larger pool in a context dependent way (sparse feature selection,⁴⁵ which we also touch on in the section on biomaterials and regenerative medicine below). This improves the ability of models to generalize well and can make interpretation easier because fewer descriptors are used in the model. Subsequently, models are usually interpreted in two ways. The first is to determine which descriptors are the most important for driving improved properties of molecules. This is called “descriptor importance” for QSAR⁴⁴ or “feature importance” for ML in general. The second, applicable to models trained on substructure-type descriptors, is to project the most important features from the model onto exemplar molecules to highlight structural features associated with more favorable activity.⁴⁶ A molecule with atoms colored according to their contribution represents a molecular “heat map.” Another important, descriptor- and model-independent method for interpreting features is to apply small perturbations to the input descriptors one at a time, while holding the other constant, and observing the effect on the modeled property (sensitivity analysis, effectively generating partial derivatives of the response with respect to the descriptors).⁴⁷ These approaches to interpretation have limitations as well.⁴⁸ It is important to recall that no statistical method can distinguish correlation from causation, and interpretations cannot always be related to a mechanism. A practical approach towards mechanistic interpretability, lateral validation,⁴⁹ is to observe trends across

related phenomena: When the choice of variables, the sign and size of their coefficients are similar across multiple QSARs, this may help mechanistic understanding and perhaps causation.

7. **ML methods.** There are many standard methods of ML in QSAR. The current wave of enthusiasm is for deep neural nets (DNN) as the ML method. Because of their relative recency and popularity across many disciplines, comparison of DNN with other popular ML approaches is presented below.

DNN methods are attractingly widespread application across many disciplines.⁵⁰ Single hidden layer neural nets were a popular ML method for developing QSAR models in the 1990's. However, neural nets have undergone a renaissance in the past decade. Algorithmic improvements, advances in hardware, use of GPUs, etc., have made DNNs practical and computationally tractable. In AI applications, such as image classification or speech recognition, DNNs have been shown to be superior to any techniques that came before. DNNs began to be applied to QSAR⁵¹ after the Merck Molecular Activity Challenge in 2012.⁵² In less than a decade we have seen an enormous growth in publications using diverse DNN architectures for modelling chemically-related properties.

To put DNNs into context for QSAR, there are many other ML methods used in QSAR modeling including k-nearest neighbors (kNN),⁵³ partial least squares (PLS)⁵⁴, support vector machines (SVM),⁵⁵ relevance vector machines, (RVM),⁵⁶ random forest (RF),⁵⁷ Gaussian processes (GP),⁵⁸ and boosting⁵⁹. In the pharmaceutical industry (in fact, in any discipline), ML and DNN methods can be compared to older methods by the following: –

1. Prediction accuracy
2. Number of sensitive and tunable hyper-parameters;
3. Need for descriptor selection
4. Length of training time
5. Length of prediction time (including uploading the model into memory);
6. Domain of applicability (determined mainly by descriptors and training set characteristics)
7. Interpretability of models.

RF has been a popular choice for QSAR modeling for many years as it can make very good predictions, has few adjustable parameters, and can be parallelized. Moreover, the degree of agreement of predictions of different agreement of RF trees⁶⁰ can help define the AD.

Boosting is also very useful because it is often one of the most accurate and fastest methods, especially with the latest implementation of Extreme (XGBoost⁶¹) and Light Gradient Boosting Machine.⁶²

The case for DNNs as a ML method would be made based on its superior predictivity. Comparison of DNNs to other ML methods like RF and XGBoost on standard industrial QSAR datasets shows a statistically significant improvement in prospective predictions as shown in studies conducted by some of the authors of this paper, and similar conclusions

have been published elsewhere.⁶³ However, in absolute terms, the improvement is less than notable. When trained on the same data sets and descriptors, DNN predictions are not different to those of other methods.⁶⁴ Thus, the squared correlation coefficient (R^2) of models generated with DNN was only 0.04 higher (on average) than those built with RF as shown in Figure 3. This is consistent with the universal Approximation Theorem discussed below.

Deep Neural Nets methods also have undesirable characteristics such as requiring more tuning of training parameters for a given training set, being computationally more demanding, taking longer to predict, and being harder to interpret.

Why are DNN models not making substantially better predictions than the other ML methods? A fundamental reason is the Universal Approximation Theorem that states that single layer neural networks (and ML methods mathematically similar) are sufficient to model any nonlinear function given sufficient data.⁶⁵ Another reason may be that any pharmaceutical data set inevitably has experimental errors that will compromise very accurate model generation. Training and test sets are also not necessarily similar, and the new field of modelability suggests that all QSAR methods are limited by the presence and size of activity cliffs.⁶⁶ For these reasons, more sophisticated and flexible methods will not necessarily provide better predictions.

It is important to remember that in the pharmaceutical industry, unlike other areas where ML is applied, the data required to build models is limited, expensive, and resource-intensive.⁶⁷ Getting marginally better predictions is not useful when the bottleneck is data paucity. However, DNNs methods do have very important advantages over most other ML methods:

1. They can straightforwardly model more than one activity at a time (multi-task models),⁶⁸ the same is true for single layer NNs with multiple output nodes⁶⁹ but not so for other ML methods. It has been claimed that on the average this produces better predictions than models of the individual activities. In practice, this effect can be quite modest, exhibiting both improvements and degradations in prediction for individual activities. It has been shown that improvement relies on the training set for the activities sharing similar compounds and features, and there being significant correlations between the activities.⁷⁰
2. Their ability to automatically generate novel chemical features (using, e.g., graph convolutional neural networks, CNNs) is particularly important.⁷¹ This mimics how images are processed on the fly (with atoms replacing pixels), as opposed to the use of pre-generated chemical descriptors. The premise is that by generating richer molecular features, more predictive models will result. In some cases, CNN has provided more accurate predictions than descriptor based DNNs.⁷¹ For example, CNN is better at predicting quantum chemical energies.⁷²
3. They provide the possibility of inverting the QSAR model (inverse QSAR), i.e. designing molecules directly from the model (so called generative models).⁷³ This is in contrast to the current QSAR practice that only goes in the direction of property prediction from structures, not from properties to predicted structures. Candidate molecules must be generated by screening large virtual libraries or by

assembling or swapping chemical fragments and predicting their properties by a QSAR model.

To summarize, it is still unclear from the ML literature whether DNNs are distinctly better at QSAR tasks than standard methods, because in most cases an exhaustive comparison has not been made. We would recommend that the method in question must always be compared to a good off-the-shelf ML method (such as RF or boosting) in the context of QSAR Best Practices.¹⁷ We would also recommend that a fairly large number of datasets (>10) should be examined in any given study. This removes the temptation to cherry-pick the results that make the method under study look better.

Another issue is the tests for DNN performance represent a low bar for success, meaning that predictivity appears better than it is in practice (an issue for the entire QSAR area). Random-split validation (which is still a literature standard) makes predictions that appear to be good because the test and training sets cover about the same chemical space, a difficult constraint as predictions outside of the model AD are likely to be poor). We recommend a time-split validation where possible, checking that the test set compounds are not too far from the model domain. Another practice in ML is to tune hyper-parameters using a validation set, where both the validation and test sets have been chosen from the same pool of compounds. In effect, this lets information about the test set to leak into the training set of the model, which makes predictions overly optimistic, and thus this practice should be avoided. The enthusiasm for DNN methods has sometimes encouraged bad practices, such as not comparing results to simpler methods (Occam's Razor) and publishing non-reproducible models, as has been reported in other areas of machine learning.⁷⁴

In our opinion the current enthusiasm for DNNs in QSAR is not yet justified by its slightly increased predictive performance, given that the methods are compute-intensive and the models very hard to interpret. However, it should not be overlooked a that their main advantage in in the generation of novel and useful features from relatively simple representations of molecules (or materials) and the potential for inverse QSAR. The development of new methods for DNN model interpretation such as Layer-Wise Relevance Propagation will also increase their advantage over traditional QSAR methods.⁷⁵ Clearly, given how fast the field is developing, it is hard to know whether DNNs will overcome current disadvantages, although the inexorable increase in computational resources available will ease some of them. On the other hand, the enthusiasm for DL methods is driving a renaissance in the use of ML in chemistry,⁷⁶ creating more opportunities.

As computational chemists, we should be actively researching other fields like data science and mathematics for advances in ML methodology. Historically, we have acquired new ML methods through serendipity, because we tend to read only the chemical literature. For example, the author of this section started applying RF to QSAR in 2003 because of a chance conversation with statisticians. We became aware of DNNs only after the Kaggle contest in 2012 and of XGBoost in 2016 because of a suggestion from a person in the IT department. However, the criteria we proposed for how DNN and ML methods should be compared, and concerns and suggestions on how best to generate dataset splits to enable robust assessment of model predictivity, have originated from our experience in QSAR

modeling. These learnings will undoubtedly be valuable for other areas of statistical data modeling. The above examples suggest that exchange of best practices and methodologies between QSAR modeling and other fields will bring advances in both. Better definitions of important general concepts such as applicability domain or model interpretability are applicable to other diverse disciplines.

QSAR in chemical safety assessment

QSAR approaches have been used extensively to model important drug properties such as ADMET. Minimizing toxicity and optimizing pharmacokinetics is critical for designing new and safe medicines; incorrect estimation of these parameters can result in undesired side effects and affect *in vivo* efficacy, leading ultimately to a failure of a drug candidate. It should be noted that almost any chemical is toxic at a sufficiently high dose, so an important characteristic of any drug is its therapeutic index, the ratio of the effective dose causing the desired therapeutic effect in 50% of research subjects (ED_{50}) to the drug dose causing adverse effect(s) in 50% of the subjects (TD_{50}). Thus, it should not be surprising that even extremely toxic compounds such as snake venom toxins are useful, at proper concentrations, as diagnostic probes, drug leads, or even as therapeutic agents.⁷⁷ Chemical toxicity is also very important for the assessment of the occupational health and environmental safety. Because toxicity is a complex multifactorial phenomenon caused by chemical effects on biological systems, it is important to understand underlying toxicity mechanisms to build mechanistically meaningful prediction models. There is a clear need to develop standardized protocols when conducting toxicity-related predictions, and the information needed for protocols to support *in silico* predictions for major toxicological endpoints of concern (e.g., carcinogenicity, acute, genetic, reproductive or developmental toxicity) across several industries and regulatory bodies has been discussed elsewhere.⁷⁸ Below, we review several key concepts that relate to issues in chemical toxicity prediction.

Adverse Outcome Pathways (AOP).

AOP is one of the key concepts of toxicity assessment. It assumes that toxicity is initiated by a molecular initiating event (MIE), which leads to an adverse outcome (AO).⁷⁹ A single AOP describes a sequence of linked events starting from MIE, going through a cascade of linked key events (KEs), and ending at an adverse health or ecotoxicological effect. The Adverse Outcome Pathway Knowledge Base is currently under active development for both health and eco-toxicology studies.⁸⁰ With knowledge of AOPs, QSAR modeling can be used to identify the potential of chemical compounds to cause a MIE and/or to lead to an adverse outcome.

Importantly, metabolites can also cause toxicity even when the precursor has low toxicity. Therefore, incorporation of information about metabolic activation can improve toxicity QSAR models.⁸¹ AOP facilitates mechanistic interpretation of models, provides a better understanding of toxicity, and allows the development of new *in vitro* tests.⁸² Currently, the development and validation of such tests is an emerging topic in predictive toxicology.

In vitro toxicity and Tox21.

Tox21⁸³ is a high-throughput toxicity evaluation initiative supported by several government agencies including US Environmental Protection Agency (EPA), National Institutes of Health (NIH), and Food and Drug Administration (FDA). Similar initiative exists in Europe under the REACH (Registration, Evaluation, Authorization, and Restriction of Chemicals) legislation. REACH encourages the use of so-called *alternative approaches* or surrogate end points to reduce animal testing. Naturally, QSAR modeling represents one of the best alternative approaches for risk assessment because it can be used both to predict *in vitro* activities of compounds and to as combine these *in vitro* results with computed molecular descriptors to improve the accuracy of models in predicting *in vivo* effects. The requirements for using QSAR models for regulatory purposes have been reviewed elsewhere.⁸⁴

Tox21 data have been used actively by the cheminformatics community to test both the prediction accuracy of QSAR models and to understand current limitations of the field. The Tox21 Data Challenge aimed to assess the ability of QSAR models to predict important *in vitro* endpoints related to chemical toxicity.⁸⁵ Participants predicted the outcomes of 12 cellular stress assays.⁸⁵ The winning team (as determined by the AUC metric) used a DNN to build multi-task models for these outcomes.⁸⁶ Model built with an Associative Neural Network⁸⁷ had similar prediction performance. The results of the Tox21 Challenge indicated that recent progress in neural networks have accelerated development of robust and predictive QSAR models for *in vitro* toxicity. The development of new types of DNN⁷⁶ has opened up new applications, allowing simpler molecular representations, such as SMILES strings or chemical graphs to be used to generate useful toxicity (and other property) models. However, these methods have generally lower prediction accuracy than ML approaches using traditional QSAR descriptors.⁸⁸ DNN methods also require substantially larger datasets to fully capitalize on their advantages⁷⁶, a problem that is rapidly abating due to explosive growth in chemical data that is driven by automation.

Tox21 data also gave rise to a number of notable comprehensive studies, such as Collaborative Estrogen Receptor (ER) Activity Prediction Project (CERAPP)⁸⁹ and Collaborative Modeling Project for Androgen Receptor (AR) Activity (CoMPARA), involving 17 and 25 international teams respectively. The resulting consensus QSAR models leveraged knowledge from the groups and were used to predict ER and AR potentials of 32,464 new chemicals.

It should be emphasized that development of new experimental techniques such as deep-sequencing RNA-Seq,⁹⁰ provides new types of data for *in vitro* assessment of toxicities that can also be used for QSAR modeling.⁹¹

***In vivo* toxicity.**

Given that adverse reactions could be caused by a multitude of factors, prediction of *in vivo* toxicity is arguably the most difficult task in QSAR modeling. The cost and ethical issues associated with direct *in vivo* toxicity assessment means that data to train models is scarce, so models are quite limited. This is clearly illustrated by the results of ToxCast Lowest

Effect Level prediction challenge.⁹² The highest prediction accuracy with the lowest RMSE of 1.08 log units was achieved using a consensus prediction of Associative Neural Network⁸⁷ models developed with several sets of descriptors.⁹² Although the organizers of the challenge have offered a set of *in vitro* measurements performed within the ToxCast project, the top-ranked model was exclusively based on the calculated descriptors and was not improved by adding *in vitro* data as descriptors.⁹² The failure of this⁹² and QSARWorld Bioavailability Challenge indicates critical importance of data curation.⁹³ Availability of more *in vivo* data, application of more complex methods such as those based on physiologically-based pharmacokinetic (PBPK) models,⁹⁴ better data curation⁹³ as well as new descriptors, which account for pharmacokinetics, should improve the model accuracy. Since *in vitro* assays in ToxCast were not predictive of such complex endpoint,⁹⁵ other methods, such as those based on systems biology, or more complex assays such as RNA-Seq used in combination with gene interaction networks, may be more successful.⁹⁶ Indeed, it was reported that combination of *in vitro* and *in silico* predictions contributed better models for a number of *in vivo* endpoints.⁹⁷

Multitask modeling: an approach that should not be overlooked.

Multitask modeling leverages information from multiple correlated properties and may provide models with higher predictive power than individual QSAR endpoint models. This is attributed to read-across and the existence of mutual information in the more complex multiple end point data sets. A recent study showed that multi-task modeling consistently improved the accuracy of models for prediction of 29 *in-vivo* endpoints using 87K chemical structures collected from the Registry of Toxic Effects of Chemical Substances (RTECS) database.⁹⁸ Importantly, authors suggested that the significantly improved toxicity predictions of multitask models should reduce the need for animal testing, prompting revisions to the current regulatory guidelines.

Structural alerts and QSAR—Identification of molecular features associated with toxicity (structural alerts) represents a tool because it can help reduce unwanted side-effects of compounds by removal of offending moieties. However, toxicity alerts generally have lower prediction accuracies compared to QSAR models.⁹⁹ It has also been suggested that a combination of alerts and QSAR models may provide improved guidance for rationally designing new compounds with reduced toxicity.⁹⁹ These combined approaches were further developed by the chemistry-wide association study (CWAS) that predicted Ames mutagenicity and an adverse drug reaction known as Stevens-Johnson Syndrome.¹⁰⁰ The identification of important chemical fragments and analysis of their co-occurrences also allows mechanistic interpretations of QSAR models without compromising their accuracy.

In summary, this section provides a brief review of a special area of QSAR modeling that deals with chemical safety. However, even in this highly specialized application there are components that can be generalized to other applications. Multi-objective modelling and optimization is one such approach that will be increasingly used in other disciplines. The ability to interpret complex statistical models for any target effect is important in many fields, especially when building models of large data sets using deep neural networks.¹⁰¹

These examples reiterate the conceptual overlap between many elements of QSAR modeling and challenges faced by other disciplines.

Multi-target profiling and polypharmacology

Since the beginning of the 20-th century, the concept of “a magic bullet” has served as the basis for drug discovery and development.¹⁰² According to this concept, a drug should be developed with the highest selectivity toward the intended target for a particular disease. Thus, classical QSAR/QSPR studies have been performed with training sets of compounds active in a single biological assay; frequently, all compounds also belong to the same chemical series.¹

The advent of high-throughput screening technologies and proliferation of diverse assays have enabled screening of a larger number of molecules in more diverse assays. Consequently, it is now generally accepted that the majority of pharmaceutical agents interact with several, sometimes many, biological targets. This often generates beneficial therapeutic activities,¹⁰² due to additive or synergistic pharmacological effects.¹⁰³ On the negative side, drugs can also interact with undesired molecular targets to causing adverse or toxic effects that often block further development. Clearly, there is a strong need to understand both the beneficial and adverse polypharmacology of ligands.¹⁰⁴

Discovery of molecules with beneficial polypharmacology could be achieved by the experimental evaluation of millions of drug-like compounds against thousands of targets. Currently, this is an unrealistic task, particularly taking into account the variability of results obtained for the same ligand-target interaction in different assays, and relatively low hit rates of experimental screens.¹⁰⁵ Thus, *in silico* prediction of biological activity profiles by (Q)SAR models is a viable alternative to these intractable experimental screens. Importantly, virtual screening approaches may be applied to millions of virtual molecules designed *in silico*. Such virtual screening greatly reduces both the number of molecules needed to be synthesized and tested, allowing pre-selection of likely hits and reduced time and cost in synthetic chemistry programs.¹⁰⁵

Multi-target profiling of compounds has led to the concept of the Biological Activity Spectrum,¹⁰⁶ defined as the set of different biological activities resulting from the compound interaction with different biological systems. It therefore represents an "intrinsic" property of the compound that depends only on its chemical structure.

Several approaches for multi-target modeling have been proposed. One of the earliest developments in this area was the computer program PASS (Prediction of Activity Spectra for Substances) reported by *Filimonov et al* almost 30 years ago.¹⁰⁷ PASS employs a uniform set of Multilevel Neighborhoods of Atoms (MNA) molecular descriptors and a Naïve Bayes classifier to model structure-activity relationships across a wide variety of biological assays. This approach allows the prediction of a wide range of biological activities at molecular, cellular, organ/tissue and organism levels. It can predict pharmacotherapeutic effects, mechanisms of action, specific toxicities, terms related to drug metabolism, gene expression, etc. The current version of PASS predicts several thousand

biological activities based on the analysis of structure-activity relationships in the training set of over one million biologically active compounds.¹⁰⁸ More recently, Gonzalez-Diaz et al.¹⁰⁹ developed the perturbation theory machine learning (PTML) methods that search for QSAR models capable of simultaneous prediction of many target properties under several experimental conditions.

Substantial amounts of relevant chemogenomics data have recently become available from PubChem, ChEMBL, and other public sources. This has catalyzed a resurgence of freely available Web-accessible tools for bioactivity predictions and continuing development of new QSAR tools and methods.

In contrast to PASS Online,¹⁰⁶ which is an open access Web-service for predicting biological activity spectra, most other tools focus on predicting putative molecular targets for compounds of interest. They use training sets extracted from publicly available data sources, different types of chemical descriptors, and prediction methods based on implementations of different chemical similarity searches. Despite some disadvantages,¹¹⁰ such approaches remain an accessible way of predicting compound activity against novel pharmacological targets lacking sufficient training data for building accurate QSAR models.¹¹¹ If the number of known ligands is sufficient for model building, some web portals provide an option to predict compound activities using conventional QSAR.

It is challenging to compare the performance of multi-target profiling tools. In contrast to single target models, there is a paucity of evaluation sets of compounds reproducibly tested for several types of biological activity. Thus, only a few comparative studies have been reported to date. For example, using data on affinity of drug-like compounds against several GPCRs, the performance of a collection of multiple target-specific k-nearest neighbors (kNN) QSAR models, PASS¹⁰⁶ and Similarity Ensemble Approach (SEA)¹¹² was compared.¹¹³ The best results were obtained with the kNN method, while PASS demonstrated a moderate predictive accuracy and SEA shown the lowest prediction power across multiple targets.

Recently, a large evaluation set including half a million compounds tested across more than 1,000 assays was constructed from ChEMBL data.¹¹⁴ The performance of several ML methods was evaluated, and again kNN generated the best results, while SEA showed the lowest predictivity. It is noteworthy that all ML methods showed relatively small differences in predictive accuracy and the advantage of the DNN was not readily apparent. This conclusion appears reasonable given that the principal purpose of DNN development was image feature recognition, i.e., similarity assessment but not prediction. Similar observations of the lack of advantage offered by DNN in cheminformatics compared to conventional ML was also made in the preceding section on modeling chemical toxicity.

As also noted in the preceding section of this paper, multi-task learning represents one of the major directions of QSAR development. A natural extension of multitarget QSAR is the analysis of ligand-target interactions in combined chemical-biological space, so called chemogenomics.¹¹⁵ Several hundred papers have been published on new methods and applications for chemogenomics (some discussed in greater detail in the following sections).

For example, Gupta-Ostermann and Bajorath reported the Structure-Activity Relationship (SAR) Matrix method, which predicts activities and allows navigation in multi-target activity spaces.¹¹⁶ March-Vila and co-workers have summarized the promise of chemogenomics applications for drug repurposing.¹¹⁷

A recently proposed proteochemometrics (PCM) approach employs relevant information from target sequences and combines it with ligand descriptors to develop models predicting ligand-receptor (class of) binding affinity. This approach is more useful than ligand-based modeling in cases when the same ligands show differential binding affinity to diverse targets. Several interesting applications of the PCM approach have been reported. For instance, this approach was used to predict ligand interactions with wild-type and mutated α -adrenoceptors where it has demonstrated superior predictivity in comparison with conventional QSAR methods.¹¹⁸ In other study, Lapins et al.¹¹⁹ applied PCM method to predict inhibition of five major drug metabolizing isoforms of cytochrome P450 (CYP1A2, CYP2C9, CYP2C19, CYP2D6, and CYP3A4) by drug-like compounds. A recent study, has also demonstrated significant advantages of PCM approach and inductive transfer of knowledge between the targets over traditional methods.¹²⁰

Careful review of the published results of PCM modeling leads to the conclusion that it may provide good estimates of ligand-target affinity in a single model by combining data from multiple assays (Figure 4). However, to achieve this goal, substantial efforts must be applied to standardization¹²¹ and curation¹²² of such data.

To conclude this section, we note that training sets used to develop conventional QSAR models do not exceed millions of entries, while the estimated size of drug-like chemical space is up to 10^{60} molecules.¹²³ We expect that with the growth of chemogenomic data and expansion of the studied chemical space, the multi-target QSAR modeling will become more common than single-target QSAR studies and that multi-target QSAR will lead to the discovery of novel medicines with much improved safety and potency profiles. Another important projection is that further development of multi-objective optimization methods will not only expand the field of polypharmacological QSAR but will also find use in many other predictive disciplines where multiple objectives need to be optimized.

QSAR-like approaches in genomics

Genomic and HTS (high throughput screening) data have rarely been subjected to QSAR analyses. Indeed, typical workflows require hit confirmation and validation prior to (Q)SAR modeling, and cheminformatics-based prioritization schemes based on individual compounds as well as scaffolds have been proposed.¹²⁴ One of the major obstacles to date remains the absence of the gene-based descriptors suitable for ML. However, high throughput driven biomedical knowledge accumulation has created an urgent need for Big Data analytics in genomics and HTS to help with the evaluation, interpretation, and integration of data, and with development of respective models.

From a life sciences perspective, the use of DNN can generate novel applications and even entirely new meaning to the field of chemical genomics by directly linking the structure of

the molecule to its effect on genes, and **by embedding these linkages in models that predict gene-mediated effects of chemicals *in vivo***. Such models require the combination of input features that characterize both small molecules (i.e., chemical descriptors) and genes (e.g., gene expression profiles) or HTS results for training. Only a few studies have been published in this area so far. For instance, it was demonstrated that gene ontology (GO) terms¹²⁵ and HTS results can be translated into input features for cheminformatics models.¹²⁶ In another such study, Sedykh et al.¹²⁷ described and implemented a workflow for using HTS data in combination with molecular descriptors to predict *in vivo* toxicity. In a related work,¹²⁸ *in vivo* rat oral toxicity was predicted by combining endpoints of 499 HTS assays (biological variables) with 548 circular Morgan descriptors (chemical variables). Notably, when used separately, biological descriptors resulted in a model with lower statistical significance than the model based on chemical descriptors.

Another example of ‘hybrid’ QSAR modelling shows how

QSAR descriptors and GO terms can be combined within a unified QSAR model capable of predicting the effect of a given molecule on a particular gene.¹²⁹ Specifically, levels of expression of 1000 ‘hallmark genes’ in six **cell lines** were predicted by DNN-classifiers, where for every molecule-gene pair in the training set, circular Morgan fingerprint values (molecular descriptors) were combined with GO terms used as gene descriptors. The resulting DNN models built with back-propagated feed-forward fully connected multi-layer perceptron (MLP) with four layers yielded good prediction accuracies (cross-validated area under the curve (AUC) values were in the 0.80-0.83 range). These results suggested that ‘hybrid’ DNN models can rather accurately associate genes and small molecules to up- or down-regulation.

Seventeen different protein- and gene- centric data sources totaling over 262.3 million data points were integrated into knowledge graph representation with typed nodes and edges, which enable the conversion of the gene-based information into descriptors suitable for ML via network-based analytical algorithms.¹³⁰ Specifically, a set of 103 genes having autophagy (ATG) associated annotations from GO terms, UniProt¹³¹ and KEGG,¹³² were used to derive ML models using the metapath approach combined with the XGBoost algorithm.¹³³ These binary ML models were trained to distinguish ATG genes from non-autophagy genes (cross-validated AUC values were in the 0.95-0.99 range). Of the top 251 predicted novel genes, 23% were associated with ATG based on literature queries, whereas 193 were not.

These case studies offer an important example of QSAR modeling evolving towards the use of more complex datasets. Synergistic use of features representing both chemical and biological properties, including gene expression profiles, GO terms and KEGG pathway associations combined with ML methods, are generating promising results. This increase in complexity is typical for many areas of research where DNN and gradient boosting methods are finding growing applicability. The improvements in model accuracy achieved by ML approaches may have been modest so far, but the prediction power of these models may increase in near future due to cross-fertilization of ideas on using ML for data modeling both in chemical datasets as well as in many other areas of science and technology. It is tempting

to speculate that DNN technology can directly screen virtual chemical libraries for compounds with bespoke, useful modulation of target genes and gene networks.

As the sources of data and sizes of datasets describing the biological properties of small molecules grow, there is also a concomitant demand for knowledge management (KM) systems, that integrate heterogeneous data into unified, predictive models and *translate data into information*.¹³⁴ For example this might allow merging of experimental bioactivity data for small numbers of molecules, 3D information from experimentally resolved structures of protein targets for these molecules, statistics of respective drug adverse event reports, and high-volume (often lower quality) data such as Genome-Wide Association Studies (GWAS) or HTS. Such large scale datasets are already assembled into knowledge graph systems, for example Pharos,¹³⁵ which supports in-depth exploration of the druggable genome.¹³⁶ Modelling such data via ML, sparse feature selection, and other advanced algorithmic approaches may lead to a better understanding of the associations between chemical structures and proteins and genes in an unbiased, objective manner They could further help identify novel gene-phenotype associations, either for diseases or for physiological phenomena such as autophagy.

QSAR in synthetic organic chemistry

The application of QSAR modeling to challenges faced by synthetic organic chemists is a recent and exciting development in predictive computational chemistry. Rapid growth in robotic platforms for drug and materials design has stimulated the development of reliable cheminformatics tools to assist with efficient synthesis of target molecules. These tools estimate synthetic accessibility of a target molecule and suggest feasible synthetic routes (Figure 5). Two of the most widely used synthesis planning strategies are forward synthesis (starting from specified building blocks) and retrosynthesis (starting from a specified target molecule). Synthetic routes usually contain multiple reaction steps for which major products and, ideally, kinetic parameters must be predicted by models. Once a given elementary reaction is selected, reaction conditions (solvent, catalyst, temperature, *etc.*) leading to a reasonable yield should be suggested by the algorithm. The above considerations can be met by a wide range of cheminformatics tools, some of which are currently used in a computer-aided synthesis design. In this section we briefly describe reaction data availability, visualization, and analysis, and summarize recent studies focused on different parts of the modeling workflow described in Figure 5.

Reaction data availability.

New modeling tools need access to large volumes of experimental reaction data stored in public and proprietary databases. In most of the recent studies, the Reaxys database (> 40 M reactions including 12.5 M one-step reactions),¹³⁷ the USP database extracted from US patents (>1.2 M reactions),¹³⁸ and the QSRR database (~10.000 reactions) have been employed. Generally, reaction data from public databases is of mixed quality. Many of the reactions are stoichiometrically unbalanced, some important data on reaction conditions are missing, and different names are used for the same catalysts or solvents.¹³⁹ However, no standards for reaction data curation have been reported so far. Ignoring the data curation step

of the modeling workflow will significantly affect the quality of the training data and models derived from them.⁹³

Reaction encoding.

Chemical reactions constitute a very complex modeling problem in cheminformatics. A reaction equation involves several different types of molecular graphs (for reactants and products) and its yield depends on numerous experimental conditions. Depending on ML method used, chemical structures can be encoded by SMILES (e.g., in sequence-to-sequence models¹⁴⁰) or by descriptor vectors, or a combined fingerprint (resulting from concatenation of descriptors of reactants and products¹⁴¹), or subtraction of descriptors of reactants from descriptors of products.¹⁴² The latter may require balanced reaction equations that, in turn, need a specific data curation step.¹⁴² Alternatively, a chemical reaction (balanced or unbalanced) can be encoded by the Condensed Graph of Reaction (CGR). This merges reactant and product structures into a single molecular graph employing both conventional chemical (single, double, etc.) and “dynamic” bonds characterizing observed transformations (e.g., single and double bond breaks, single-to double bond conversion, etc.).¹⁴³ CGR can be considered a *pseudomolecule* to which any cheminformatics approaches can be applied. In particular, fragment descriptors or fingerprints can easily be generated for CGR.¹⁴⁴ Solvent can be encoded by a set of physicochemical parameters which can be concatenated with the structural descriptors.

Visualization and analysis of reaction space.

Both graph-based and vector-based approaches have been used to visualize the chemical space of reactions. In graph-based approaches, chemical reactions and individual molecules (reactants and products) are represented as nodes of a large bipartite graph¹⁴⁵ used to optimize synthetic pathways. In the vector-based case, a chemical reaction is defined as a vector in multidimensional space defined by descriptors. Dimensionality reduction is required to generate a two-dimensional map describing the data distribution. This approach was pioneered by Gasteiger et al.¹⁴⁶ who generated Self-Organized Maps (SOM) that clustered different classes of reactions effectively. Generative Topographic Mapping (GTM) approaches have recently been used to visualize large sets of S_N2 , cycloaddition, and tautomerization reactions. Unlike SOM and many other dimensionality reduction methods, GTM can be used to predict properties of new reactions projected on the map. As a predictive tool, GTM performs similarly to conventional ML methods like SVM.

Planning organic synthesis using prediction of reaction products and retrosynthetic analysis.

The general aim of synthesis planning is to identify a series of feasible reaction steps leading to a target compound from available starting materials. Retrosynthetic methodology, invented by Corey,¹⁴⁷ is a real challenge because the search for precursors of a product generates a combinatorial explosion of possible reaction routes. Cheminformatics tools can help select the most feasible series of single-step reactions. The current trend in this field is to train DL models on large sets of reactions to predict probabilities of different retrosynthetic transformations. It was shown that using Monte Carlo tree searches and symbolic AI methods, it is possible to identify feasible reaction pathways.¹⁴⁸

Prediction of reaction outcomes allows one to prioritize retrosynthetic suggestions. A cheminformatics tool should predict the products of a given set of reactants under given conditions. Consideration of multistage chemical transformations and competitive reactions will significantly complicate this problem. Current trends in the modeling of reaction outcomes focus on processing large reaction databases with DL models to predict the probabilities of competitive chemical processes.¹⁴⁹ The latter can be used directly for reaction outcome predictions. The ReactionPredictor tool¹⁵⁰ is of particular interest because it forecasts the output of complex chemical reaction by combining mechanistic considerations with ML. This approach enumerates possible interactions and then ranks them using a pseudomolecular orbital approach.

Two orthogonal methodologies, template-based and template-free, can be applied to retrosynthesis and outcome prediction. Template-based methods rely on user-established sets of transformation rules, either suggested by expert-chemists or extracted automatically from reaction databases, the feasibility of which is assessed by the model. This concept is employed in most retrosynthetic tools, including the popular CHEMATICA program,¹⁵¹ which integrates more than 10,000 empirical transformation rules.

Alternatively, in template-free approaches transformations between the reactants and the products of chemical reactions are deduced directly from their structures. This allows one to automatically enlarge the list of transformation rules as soon as new data are available. This methodology has become more popular in recent years. For instance, Coley et al.¹⁵² suggested using a graph-convolutional neural network and a global attention mechanism, followed by the application of rules to reaction product predictions and retrosynthetic analysis. Another template-free approach employs natural language processing methods, namely 'sequence-to-sequence' models. These use recurrent neural networks (RNN), commonly applied to translation of texts between languages. When applied to chemical reactions, SMILES strings of reactants and products constitute the language. This methodology was applied to model reaction products and for retrosynthetic reaction route prediction, which provided similar performance (ca 37% for top-1) to rule-based systems (35%).¹⁴⁰ A use of an advanced Transformer architecture, which was initially used for English-to-German translation, boosted the accuracy of predictions to about 43%.¹⁵³ This result indicates that retrosynthesis predictions can be significantly improved by algorithms originally developed for very different purposes.

Forward synthesis planning.

One of the most impressive approaches to forward synthesis planning has been implemented in the DOGS program.¹⁵⁴ This algorithm applies 58 well-established chemical transformation rules to a set of 25144 readily available synthetic blocks from the Sigma-Aldrich catalog. New molecules are grown in a stepwise procedure, each step consisting of complete enumeration of all possible solutions followed by selection of top scoring intermediate products to subsequent growing steps. The quality of designed products is assessed using pairwise similarity to a target molecule. Thus, DOGS can usefully suggest a synthetic plan not only for the target molecule but also for its close analogs

Assessment of synthetic accessibility.

Synthetic accessibility (or the opposite, synthetic complexity) is a scoring metric used to prioritize virtual compounds for synthesis. It is often used as an important filter for screening virtual libraries and in *de novo* design studies. Among scores developed so far¹⁵⁵ the most popular is SA score.¹⁵⁶ It is calculated using contributions from fragment occurrences in PubChem compounds and a complexity penalty based on the number of chiral centers, rings, macrocyclic fragments, and the total number of atoms. Recently, Coley et al.¹⁵⁷ suggested the Synthetic Complexity Score (SCS) which relies on a neural network trained on 22 million reaction pairs from the Reaxys database.

Prediction of kinetic and thermodynamic characteristics.

The logarithm of the reaction rate constant ($\log k$) is a common endpoint in QSAR modeling, first used more than 70 years ago.¹⁵⁸ Currently, quantitative structure-reactivity relationship (QSRR) modeling is performed on large and diverse datasets that account for solvent effects and temperature for many types of chemical reactions using NN approaches. In these models, descriptors computed for the reactants are concatenated with solvent and temperature descriptors. This technology must know the order of reactants in the reaction equation, making the development of an automatized QSRR workflow problematic. This problem can be solved using Condensed Graphs of Reaction (CGRs) that combine the reactant and product information. Fragment descriptors generated for CGRs were concatenated with solvent and temperature descriptors and used to train $\log k$ models for bimolecular nucleophilic substitution,¹⁵⁹ bimolecular elimination, and different types of cycloaddition.¹⁶⁰ Similar approaches were used to develop predictive models for the equilibrium constants of tautomerization reactions.¹⁶¹

Prediction of optimal reaction conditions.

Since the reactivity of chemicals is largely determined by the reaction conditions, their theoretical assessment is of particular importance (especially for automated robotic synthesis). Several approaches to reaction conditions modeling have been reported. For example, Marcou et al.¹⁶² used CGR-derived fragment descriptors to train SVM, RF, and Naive Bayes classification models to predict optimal solvents and catalysts for the Michael reaction. Gao et al.¹⁶³ reported NN-based models trained on ~10 million reactions from Reaxys that identify appropriate catalysts, solvents, reagents, and temperatures for a specified reactions. A 70% match with experimental conditions was found within the top-10 predictions. Lin et al.¹³⁹ used the heuristic that similar reactions proceed under similar conditions to predict optimal reaction conditions. They used a simple similarity search of reaction databases with recorded conditions,¹³⁹ especially effective with the CGR technology.¹⁶⁴ The value of this approach has been demonstrated by protective group deprotection reactions. Models trained on 142,111 catalytic hydrogenation reactions demonstrated high accuracy (ca. 90%) for predicting optimal experimental conditions.

In summary, the CGR technology can efficiently model optimal reaction conditions. One employs similarity searching of reaction databases to construct QSRR models, with reaction conditions as endpoints. Studies summarized in this section provide compelling examples of the impact of QSAR modeling on one of the historically most empirical areas of natural

science, synthetic organic chemistry. The development of both retrosynthetic and forward synthesis prediction models, based on the analysis of an immense amount of accumulated data, represents one of the most important frontiers in modern science. It is essential for chemists to understand and begin applying these emerging approaches. When coupled with robotic synthesis methods, these synthesis prediction models are poised to transform organic chemistry as we know it and open the door to autonomous chemical synthesis systems in the future.

Closed-loop discovery and automation

Traditional serial molecular and materials discovery processes in laboratory have arguably reached a plateau. The costs of discovering materials and drug candidates remain high and the discovery and translation time is still long. Three decades ago, combinatorial chemistry (also known as high-throughput experimentation, HTE) promised to reinvigorate the discovery pipeline by carrying out synthesis and experimentation rapidly, in parallel using automation. HTE led to important discoveries (such as novel polymers) and, indeed, has accelerated the discovery pipeline. However, the avalanche of new drug leads that was anticipated did not occur. More recently, DNA-encoded chemical libraries have made possible synthesis and testing of millions of compounds¹⁶⁵ and many big pharma companies have embraced this approach.

Furthermore, there is a growing realization that experimentation can be analyzed in terms of information theory. Questions like *What is the amount of information that an experiment contains? What is the next best experiment to carry out?* can be answered by modern Bayesian methods. This thinking has led to the revival of methods for developing *closed-loop or autonomous* approaches. By closed-loop we mean that the experimental system is designed using an information-theoretical approach, and the experimentation and assays are carried out in an automated way. By using AI or evolutionary algorithms to make decisions on what compounds to synthesize in the next cycle, in principle, an *autonomous* system can be developed. The term “self-driving laboratory” has been also coined to describe this type of experimental setting.¹⁶⁶ Clearly, a self-driving closed loop laboratory is fundamentally different from existing HTE. The closed-loop approach, designed to provide rapid iterations using autonomous decision making, seeks to *minimize* the number of experiments required to reach a specified goal (e.g., target molecule(s)). It does not need to create large libraries, rather employs agile experimental infrastructure, and statistics and ML to build QSAR-like models to predict the target properties for every element of the self-driving laboratory.¹⁶⁷

Bayesian methods show promise for making closed-loop decisions. Based on prior assumptions about the nature of the experimental observations, they can propose the optimal next experiment to conduct. PHOENICS,¹⁶⁸ for example, employs Bayesian Neural Networks and a kernel density estimate approximation to balance exploration vs. exploitation. Human interpretability is also an important factor in these systems. The algorithm chooses a set of experimental conditions to be generated by robot synthesizers. It is not sufficient to understand *what* the system generates but *how?* Interpretability is clearly very important for modern ML research. To aid interpretability, researchers have used hierarchical optimization approaches that operate on one or more variables. In multifactorial

systems it is often necessary to understand the pareto-optimal regions of the problem space. A mathematical function called CHIMERA was recently introduced to address these problems;¹⁶⁹ it can be used with any optimizer, such as PHOENICS.

Such systems require an operating system that is open-source and capable of controlling experimental equipment, storing data in databases, coupling with optimization approaches, and interacting with researchers. A “Cortana” or “Alexa” digital assistant for scientists that is connected to the closed-loop system could accelerate adoption and innovation. Efforts such as ChemOS can help rally developers to achieve this vision.¹⁷⁰

One of the promising applications of closed-loop discovery is in the materials space. A recent review summarized the state-of-the-art and challenges in this field.¹⁷¹ Examples of the application of AI to materials discovery are described in this review, as well as in following sections of this paper. One such example is the design of blue emitters for organic light-emitting diode devices accomplished by virtual screening of half a million molecules .¹⁷² This approach led the successful discovery of three lead candidate compounds with state-of-the-art performance,¹⁷² exemplifying the promise of closed loop discovery. The three good candidates required the synthesis of only ~40 materials. In autonomous systems, experimentation becomes the bottle neck in the accelerated discovery process. This can be overcome by technological developments – creation of self-driving, closed loop robotic laboratories controlled by AI, as discussed in a recent perspective.¹⁷³

Evolutionary algorithms can also be used to generate closed loop, autonomous molecule and materials discovery system. Their application to drug discovery and optimization, and materials discovery have been reviewed recently.¹⁷⁴ ML-based QSAR can be used to model the fitness landscape of materials experiments, which can substitute for downstream experiments, improving efficiency and speed.

In summary, AI methods and models that optimally instruct every step of robotic synthesis (including the choice of both reagents and reaction conditions) represents a landmark in the extension of QSAR methods toward dramatically more efficient chemical synthesis.

Machine learning approaches in quantum chemistry

Computational chemists, physicists, and biologists commonly employ molecular potentials to evaluate energies and forces. These are used to search for novel drug compounds and materials. Hence, a faster but still accurate computational method for evaluating molecular potentials is a very important development. Potential applications include calculating the free energy of protein-ligand binding via molecular dynamics simulations, and the simulation of deformation dynamics in materials.

The potential energies and forces provided by molecular potentials are obtained traditionally by quantum mechanical (QM) calculations or classical physics-based force fields (FF). QM methods solve the Schrödinger equation and are the most accurate methods for describing atomistic systems. The high computation cost of QM and long-time scales relative to experiment has limited studies of larger, realistic atomistic systems. Hence, novel robust approaches approximating QM methods without any loss in accuracy are required for

continued scientific progress. Force fields are computationally efficient, allowing the simulation of up to millions of atoms, but they require explicit parametrization of classical bonding, angle, torsion, and possibly higher-order terms. The correct parametrization of force fields can be tedious and cumbersome. Further, parametrization for one atomistic system may not be transferable to new systems.

Recent breakthroughs in the development of ML methods in chemistry¹⁷⁵ have produced general purpose models that predict potential energies and other molecular properties accurately for a broad class of chemical systems. General purpose models promise to make ML a viable alternative to classical empirical potentials (EP) and force fields since EPs are known to have many weaknesses, such as poor description of the underlying physics, lack of transferability, and are hard to systematically improve their accuracy.

Molecular representations.

To develop a useful and efficient ML-based property predictor, the most critical issue is how to represent the system in question to a ML method. These representations (descriptors) consist of some numerical representation of a molecule or a system of atoms. There are a wide range of published descriptors such as the Coulomb matrix¹⁷⁶, or its recent Bag of Bonds (BoB)¹⁷⁷ extension. Other popular choices include descriptors that represent molecular graphs,¹⁷⁸ bonds and angles,¹⁷⁹ many body expansions,¹⁸⁰ the atomistic local chemical environment,¹⁸¹ and end to end models that learn the best description of the system given minimal neighborhood information.¹⁸² Many of these techniques have been successfully applied to either molecules or materials.

Some recent descriptors like MBTR (Many-Body Tensor Representation) and SOAP (Smooth Overlap of Atomic Positions)¹⁸³ can describe both finite- and periodic systems. MBTR is derived from the Coulomb matrix, BoB, and many-body expansion. SOAP kernel represents the local density of atoms within the environment as a sum of Gaussian functions centered on each of the neighbors of the central atom. It essentially defines the similarity between two neighboring environments and uses it as a descriptor for ML models.¹⁸⁴

Local atomic environment vectors (AEV) are another widely used molecular representation. AEV explicitly includes all pairwise combination of elements, which means that the size of the input layer of a ML model grows as $\mathcal{O}(N^2)$ with the number of included chemical elements. Therefore, models can only be trained for a relatively small number of chemical elements. Adding new elements requires retraining the ML model again from scratch.

Recently, alternative weighting functions (wACSFs),¹⁸⁵ circumventing the above issue, have been proposed. Though this is a simple re-parametrization, the number of required symmetry functions becomes independent of the actual number of elements present in the system, leading to more compact descriptors. This alternative solution to the growth problem was introduced with the Deep Tensor Neural Network (DTNN)¹⁸⁶ and Atom-in-Molecule Neural Network (AIMNet). These constitute learnable vectors of atomic features that are used to embed atomic symmetry functions to make a unified representation of each atom's chemical environment. DTNN was subsequently refined to create the SchNet architecture¹⁸²

specifically designed to model atomistic systems using continuous-filter convolutional layers.

Neural Network potentials.

A ML approach applicable to chemical systems containing large numbers of atoms, originally proposed by Behler and Parrinello (BP method) in 2007, used high-dimensional neural network potentials (NNP).¹⁸⁷ As in many conventional empirical potentials, the potential energy E is the sum of local atomic energies of all atoms in the system. Since this seminal publication, a substantial number of articles and reviews have been published on the use of NNPs for bulk chemical systems (e.g., bulk silicon or water) or for describing single molecule potential energy surfaces and reaction coordinates.¹⁸⁸

Recently, Smith et al. introduced the first NNP designed for organic molecules, ANI-1.¹⁸⁹ It is applicable to molecular systems well outside its training set. The ANI-1 potential was trained on a dataset of small organic molecules of up to 8-heavy atoms (while sampling both conformational and configurational space). Furthermore, ANI-1 demonstrated its applicability to much larger systems, up to 70 atoms, including known drugs and molecules randomly selected from the GDB-11¹⁹⁰ database and containing up to 10 heavy atoms. It predicted DFT energies of the test set molecules with up to 10 heavy atoms very well, with the resulting RMSE values below 0.57 kcal/mol.

Many techniques for improving the accuracy and transferability of general-purpose ML potentials have been employed. Among these, active learning methods, already proven successful in conventional QSAR modeling, have been especially popular.¹⁹¹ Active learning methods provide a consistent and automated improvement in accuracy and transferability and have contributed greatly to the success of general-purpose models. An active learning algorithm decides what new QM calculations should be performed then adds the new data to the training set. Allowing the ML algorithm to drive sampling improves the transferability of an ML potential greatly. Further, transfer learning methods allow the training of accurate ML potentials by combining multiple QM approximations.

One fundamental limitation of BP-type models is the inability to pass information between atoms at larger distances. Several neural network architectures have been proposed to address this limitation. The HIP-NN (Hierarchically Interacting Particle Neural Network) approach breaks molecules down into feature representations and uses a number for each atom and the pairwise distances between atoms. On-site layers encode information specific to each atom and interaction layers allow sharing of information between nearby atoms. The total energy is built hierarchically from those interactions.

Another architecture, SchNet, encompasses atom embeddings, interaction refinements, and atom-wise energy contributions. At each layer, the atomistic system is represented on atom-wise basis and is refined by continuous filter convolutions with filter-generating networks.¹⁹²

In the AIMNet implementation, the solution to the short-range problem is inspired by mean field theory (MFT). The main idea of MFT is to replace all interactions with any one atom

with an average or effective interaction, sometimes called a molecular field. This reduces any multi-body problem into an effective one-body problem.

Datasets.

As previously stated, one of the most important aspects of building a model in chemistry is the choice of the training dataset. Various datasets of organic and materials systems for training ML models have been developed over the last decades. Two of the most popular organic molecule benchmark sets are the QM7¹⁷⁶ and QM9¹⁹³ collections. The QM7 benchmark was developed by subsampling the GDB-13¹⁹⁴ database of small molecules. QM7 contains 7165 energy-minimized molecules consisting of up to 7 heavy atoms and several properties computed with density functional theory (DFT). This benchmark is difficult to model by ML because of its relatively small size. Initial mean absolute errors (MAE) ,using the coulomb matrix representation,¹⁷⁶ were around 10 kcal/mol.

The ANI-1 dataset includes organic molecules with a large number of non-equilibrium DFT total energy calculations . It includes ~24M conformations for 57,462 molecules from the GDB database, with the total energy values computed for each conformation. This dataset samples both chemical and conformational degrees of freedom at the same time and thus provides 100x more data. Therefore, we expect that this dataset will become a new standard for comparing the ability of current and future ML methods to improve on the best model accuracy (1 kcal/mol) achieved for the QM9 benchmark. More importantly, this data source is a foundation for development of future general-purpose machine-learned approaches.

The COMP6 benchmark dataset¹⁹¹ was developed to validate the transferability of ML potentials. COMP6 is a benchmark suite containing five rigorous benchmarks that cover broad regions of organic and bio-chemical space of isolated molecules and a sixth built from the existing S66x8¹⁹⁵ noncovalent and intermolecular interactions data.¹⁹¹ Properties are calculated using the ω B97x/6-31G(d) basis set, however, it could be recomputed using any desired quantum level of theory.

Advanced approaches

In addition to active learning, there are other ML techniques that aim to reduce training data requirements. Some ML-based methods (such as NN) can take advantage of information from multiple sources. The key concept is to train a model using a large dataset of medium accuracy, then retrain the model with a smaller, more accurate and difficult to obtain data set. This process called transfer learning (TL) relies on the assumption that less accurate data sets contains some information that makes it easier to learn models for the smaller datasets of higher accuracy data.

For example, TL could be performed by taking a DL model that was pretrained to medium-fidelity DFT, holding some number of parameters in the model constant, then retraining the remaining parameters using a much smaller, higher accuracy CCSD(T)/CBS dataset. Such methodology resulted in the development of the ANI-1ccx potential, which represents an attractive alternative to DFT and standard force fields for conformational searches, molecular dynamics, and the calculation of reaction energies. The computed reaction energy

values demonstrated that the transfer learning-based ANI-1ccx method outperforms DFT on test cases, especially those where DFT fails to capture reaction thermochemistry.

In many systems, multiple data modalities can be used to describe the same process. One such physical system is the human brain, which provides more reliable information processing based on multimodal information.¹⁹⁶ Many ML related fields of research have successfully applied multimodal ML model training.

In chemistry, molecules, often represented by structural descriptors, can also be described by accompanying properties (dipole moments, partial atomic charges) and even electron densities. Using multimodal information as inputs has been an actively developing field in recent years.¹⁹⁷ This boost is caused by the use of additional information that captures the implicit mapping between the learnable endpoints. We discussed the advantages of multi-objective models over traditional single task approaches in the sections on Chemical safety prediction and Multi-target profiling above. Here we show that the same approaches are equally useful for developing ML models of QM results.

In the previous sections we have commented on the ongoing revolution in organic chemistry brought about by advances in computational (retro)synthetic approaches and robotic chemistry. Similarly, the use of ML approaches in quantum chemistry constitutes another recent paradigm shift. These rapidly emerging approaches dramatically change current limits of the size and complexity of molecular systems accessible to QM-level structure and property calculations.

Materials informatics

Machine learning methods dependent on large experimental and computational databases, are becoming ubiquitous tools for materials development,¹⁹⁸ extending their traditional use for organic molecules. Materials science is a very large field and space constraints permit discussion of only a small set of important *questions and answers* described below.

Which materials are missing?

This has been a perennial question,¹⁹⁹ but several recent studies have attempted to address this. For instance, Hautier *et al.*²⁰⁰ used experimental data to create a probabilistic framework for ionic substitution capable of dealing with sparse spaces (quaternary configurations). ML has also been used to tackle amorphous systems. For example, Perim *et al.*²⁰¹ identified an energy spectral descriptor for de novo prediction of metallic glasses and used it to quantify the classification probability of mixtures. ML and atomic features (descriptors) were also used to identify regions of compositions prone to glass formation and demonstrated surprising accuracy.²⁰²

Descriptors, the Holy Grail of optimization: where can we find them?

While the great importance of descriptors has been established,¹⁹⁸ these parameters are often defined *deus-ex-machina* out of intuition. Attempts have been made to develop interpretable parameterizations with ML. Thus, Ghiringhelli *et al.*²⁰³ proposed compressive sensing to discover functional forms and tested stability rules for binary semiconductors.

Isayev *et al.*²⁰⁴ introduced universal fragment descriptors for predicting properties of inorganic crystal and developed electronic density of states and band structure fingerprints that cluster many high temperature superconductors (materials cartography). Recently, Stanev *et al.*²⁰⁵ identified 30+ non-cuprate and non-iron-based oxides, potential new superconductors, using RF.

Can enthalpies (and other properties) be predicted?

The correct calculation of enthalpies and other properties is important for *ab-initio* computational materials design.²⁰⁶ Much progress has been made since the original principal components analysis of alloy thermodynamics reported by Curtarolo *et al.*²⁰⁷ Rupp *et al.*²⁰⁸ used kernel ridge regression for modeling molecular atomization energies with mean absolute error of ~10 kcal/mol. In a related study, De *et al.*¹⁸⁴ used the smooth overlap of atomic positions (SOAPs) to introduce a very useful descriptor for comparing structures: the “alchemical similarity” for molecular and periodic structures. Gaussian process regression (GPR) was used to generate very accurate Gaussian atomic potentials (GAP) and then to train a SOAP–GAP model within a ML framework (GPR) that achieved a 99% accurate atomic-scale properties for Si surface reconstruction, stability of molecules, and protein ligands.²⁰⁹ Pilia *et al.*²¹⁰ tackled melting temperatures of the octet subset of *AB* solids and band gaps of double perovskites. De Jong *et al.*²¹¹ used statistical learning to study elastic moduli of inorganic crystals, and with many other relevant studies.

What material properties can we predict?

Thermoelectrics.—A lot of work has been performed for computational predictions of thermoelectric systems following the seminal paper of Madsen who proposed an automatic search for new thermoelectric materials leading to LiZnSb.²¹² Legrain *et al.*²¹³ developed a ML descriptor-based framework (random forests and nonlinear support vector machines) and found that chemical composition alone can reasonably predict vibrational free energies. In the work of Carrete *et al.*,²¹⁴ authors used classification trees to address nano-grained half-Heuslers thermoelectrics.

Magnets.—In Sanvito *et al.*,²¹⁵ the ideal latent heat curvature introduced in Yong *et al.*²¹⁶ was calculated for all the Heusler configurations of the AFLOW repository. This was performed with the cloud phase diagram calculator by Oses *et al.*,²¹⁷ leading to the discovery of two magnets Co₂MnTi and Mn₂PtPd, the first ever discovered by computational means. Körner *et al.* performed a ML high-throughput-screening of intermetallic ThMn₁₂-type phases and rare-earth-lean systems with YNi₉In₂-type.²¹⁸ Möller *et al.*²¹⁹ built kernel-based ML models to optimize chemical compositions for permanent magnets.

Light conversion and emission.—To overcome input constraints of common ML pipelines, Duvenaud *et al.*²²⁰ developed a convolutional neural network operating directly on graphs (representing molecules of arbitrary size/shape), demonstrating enhanced predictive performance over traditional fingerprinting for solubility, drug efficacy, and organic photovoltaic efficiency datasets. Gómez-Bombarelli *et al.* integrated neural networks as part of a larger computational discovery pipeline to prioritize molecules for quantum simulations.

¹⁷² This led to the discovery of molecular organic light-emitting diodes with external quantum efficiencies as large as 22%.

High-entropy systems.—High-entropy materials continue to attract research interest due to their remarkable properties, and several semi-empirical methods have been proposed to predict their existence.²²¹ Most approaches use descriptors with parameters fitted to the limited experimental data. Modeling phase diagrams with CALPHAD also suffers from insufficient experimental knowledge.²²¹ There was a recent attempt by Lederer *et al.*²²² to parameterize the miscibility-gap and solid-solution boundary lines with ab-initio calculations and statistical modeling. Eventually, such analysis might mature into effective ab-initio descriptor-based characterization.

Other notable applications.—Fernandez *et al.* proposed an innovative QSPR model to recognize efficient metal organic frameworks for CO₂ capture. Emery *et al.*²²³ performed a descriptor based combinatorial analysis of perovskites for thermochemical water splitting applications.

2D materials.—Single or multiple layers of the same or different 2D materials have exciting new electrical, optical, heat transfer, and lubrication properties. Recently layers of graphene have exhibited superconductivity.²²⁴ ML methods have been used to predict the interlayer distance, band gap, thermodynamic properties and superlubricity properties of hybrid 2D materials.²²⁵

Welcoming new challenges!

Materials science properties, based on fundamental principles, are intrinsically suitable for modeling by machine learning. Success in ML approaches is a driver for the discovery and/or optimization of new materials and/or phenomena. This section has given a short — unavoidably incomplete — snapshot of the current state of the art.

What do our colleagues say about future frontiers?—Jain *et al.*²²⁶ identified challenges as follows: (i) streamlining the use of large data resources (even with rational APIs, large databases remain difficult to interrogate, especially when mixing data from different repositories); (ii) developing descriptors for crystalline, periodic solids; and (iii) balancing interpretability (physical meaning) of descriptors versus accuracy of models. The latter represents a well-known challenge resolved in cheminformatics about a decade ago.²²⁷ Butler *et al.*¹⁷⁵ added the following extra challenges to this list: (iv) dealing with smaller datasets (of critical importance especially for the experimental world); (v) quantum learning (to enhance calculation speed); and (vi) establishing new principles (not only data, but also laws, somewhat similar to Jain's point about balancing interpretability and accuracy).

What do we say about future frontiers?—There is no need to add further elements to the philosophical discussion of ML/AI's future. We should not underestimate the critical issues of the following additional challenges: (vii) dealing with the disordered/amorphous systems (e.g., it is not a coincidence that the field of high-entropy alloys is still lacking a compelling ML work); (viii) sustainability and organization of big-data in terms of

computational infrastructure, standardization of data-entries and prototypes, development of materials database languages, e.g., AFLUX,²²⁸ (ix) further exploration of web-, cloud-, and frameworks-directions, and the *last but the most important point* (x) unless ML can generate new useful materials faster than experiments alone, materials scientists' interest in ML will dissipate quickly.

To conclude this section, we highlight the clear similarity between materials informatics with the traditional workflow of QSAR modeling (see Figures 1 and 6). As with cheminformatics, the starting point of materials informatics is the accumulation of large datasets of materials with experimental or computational properties. The need for developing novel materials descriptors and their use in building property prediction models using ML techniques follows. Finally, current challenges outlined in the concluding part of this section parallel many of those facing traditional QSAR modeling of bioactive compounds. Thus, materials informatics (and a closely related field of nanomaterials informatics described in the next section) represents a prime example of a new discipline, whose development was enabled and immensely catalyzed by the experience and approaches developed in QSAR.

Nanomaterials informatics

Nanotechnology is another field for which cheminformatics is becoming a key tool, especially for the quantification of diverse properties of nanomaterials and nanostructure-property modeling. Development of modern AI algorithms has stimulated an increased interest in Quantitative Nanostructure – Activity Relationships (QNAR)²²⁹ also known as nano-QSAR. Like traditional QSAR, QNAR models are based on the assumption that similar nanomaterials will induce similar biological effects. However, unlike QSAR, nanomaterials (and materials in general) are more complex than single drug molecules, as they are less well defined and feature distributions of sizes, shapes, etc.

QNAR models rely on an ensemble of molecular descriptors that encode constitutional, topological, or geometrical characteristics of a given set of nanomaterials. These descriptors are derived directly from the structures of the nanomaterials using bespoke software. Moreover, experimentally determined properties (e.g., elemental composition, zeta potential, size distribution, shape) can also be appended to the computed descriptors to boost the prediction performances of QNAR models. This is analogous to the use of experimental HTS results as descriptors to model biological endpoints for drug candidates described in prior sections. QNAR models establish quantitative relationships between those experimental and computed descriptors and specified biological endpoints using ML techniques.

Importantly, QNAR models are developed using the same workflow, validation procedures, statistical criteria, and key steps as those of classical QSAR models for small molecules. However, the high structural diversity and complexity of nanomaterials typically lead to specific challenges,²³⁰ especially when it comes to the choice of molecular descriptors. Two types of representations are clearly emerging from the literature – studies in which the whole nanoparticle is characterized computationally, experimentally, or both or when such

characterization is applied to the surface chemistry of the nanoparticle (especially, organic decorators) only. Naturally, the choice of descriptors and the associated software is different for these two types of QNAR modeling. For the second type of study the QNAR model is similar to a traditional QSAR model, trained using descriptors for surface chemistry, to predict biological activity of the nanomaterials. Another challenge of QNAR modeling, similar to materials informatics is the relatively small size of the datasets currently available in the public domain. This leads to lower prediction accuracy and smaller applicability domains for QNAR models compared to those of QSAR models trained on large organic molecule data sets. To mitigate this limitation, read-across techniques are increasingly used to estimate the properties of nanomaterials.²³¹

Assessing the environmental impact of engineered nanomaterials (ENMs, see Figure 7) requires data on their physicochemical and bioactivity properties, as well as bioaccumulation. After data collection and validation, ML approaches can be used to generate models correlating values of ENMs descriptors (e.g., structural, physicochemical, and bioaccumulation-related) and specific toxicity outcomes associated with biological mechanisms of action under various exposure scenarios.

The importance of data on ENMs structure and properties

Like other area of materials science, nanotechnology has generated various datasets of physicochemical properties, environmental fate and transport parameters, and bioactivity of nanomaterials.²³² They contain both literature curated and raw data from various experimental investigations, useful for QNAR modeling. For example, the OCHEM database²³² contains experimental data on ENMs and provision for generating descriptors for model building, NanoMiner²³³ contains data (including omics data) on 634 types of ENMs. The NM-Biological Interactions Knowledge base contains over 200 toxicological evaluations for embryonic zebrafish exposed to metal and metal-oxide ENMs. NanoDatabank²³⁴ has raw data for over 1000 different nanomaterials and associated characterization and toxicity data.

Early nanoinformatics efforts were focused on organizing data into *structured* datasets (i.e., with fixed fields or records).²³⁵ However, there is growing recognition that significant data are available as *unstructured* datasets (i.e., with no predefined fixed fields or records), often are scattered across multiple literature and online sources. Thus, significant recent efforts have been devoted to the development of public databases, meta data, and data management systems for nanomaterials. These efforts included incorporation and integration of information from multiple sources, addressing data security, effective data sharing, intelligent data queries, and data integration.²³⁶ The joint EU-US Nanoinformatics Roadmap 2030²³² has stressed the need for guidelines concerning the development of nanoinformatics datasets that are structured, have controlled ontology for ENMs properties and bioactivity, and interoperability with other databases and modeling tools. Raw data (free from pre-processing by data curators) that can be curated and analyzed in a context-dependent way are most useful for QNAR development.

Substantial amounts of experimental data on the toxicity of ENMs have been generated, primarily in various cell lines such as, macrophages, pancreatic and other human cells and

bacteria. There are still limited studies with simple organisms like zebrafish and even fewer on higher animals. Toxicity data include experimental results across multiple assays and cell lines/types with ENMs having different surface modifications and core compositions. There are different levels of confidence and consistency across the toxicological studies. Currently, efforts to derive generalized toxicity models based on ENMs characteristics have been based on datasets from single studies rather than integrated from the collective body of published data.²³⁷ Clearly, to develop predictive nano-SAR models of ENMs toxicity, it is useful to identify critical biological pathways that can lead to adverse outcomes.²³⁸ Understanding relationships between the structural and physicochemical properties of ENMs and the biological responses and correlation between such responses can be very useful for deriving causal relationships. Although QNAR models provide valuable insight on ENMs toxicity, they generally cannot provide direct mechanistic interpretation that can be validated and tracked back directly to experimental data. However, as with most other QSAR models, ENMs toxicity models can be very useful in the absence of mechanistic information or interpretation.

Clearly, to generate the most robust and predictive ENMs toxicity models, the quality of data is paramount. These models can then elucidate the relevance and significance of ENMs properties such as structure, surface chemistry, shape and other physicochemical parameters with respect to their biological properties. Experimental conditions can also be employed as independent variables when modeling toxicity. Several literature studies have identified causal relationships between the biological outcomes and important ENMs properties.²³⁹

QNAR modeling

Several seminal publications pioneered the field of QNAR modeling. Puzyn et al.²⁴⁰ built the first nano-QSAR model based on ensemble learning regression methods and CDK descriptors to predict the cytotoxicity of 17 unique metal oxide nanoparticles. Fourches et al.²⁴¹ introduced the concept of QNAR modeling with a set of 109 functionalized CLIO nanoparticles and their Paca2 cell uptake. This study has been repeated and successfully reproduced several times by other research groups.²⁴² For instance, different series of metal oxides were also modeled using the OCHEM webserver to generate reliable QNAR models.²⁴³ Drug delivery properties of nanocarriers could be successfully predicted by QNAR models as well.²⁴⁴

Important nanomaterials, carbon nanotubes, have had their biological effects extensively modelled by QNAR. For instance, Trinh et al.²⁴⁵ used a combination of computed and experimental descriptors, encoded as quasi-SMILES, to build QNAR models that could accurately estimate the cytotoxicity of carbon nanotubes in human lung cells. Fourches et al.²⁴⁶ developed a series of QNAR models for 83 functionalized CNTs tested *in vitro* for protein binding and toxicity. These models reached prediction accuracies up to 74% for external test set toxicity estimates, and protein-binding classification models achieved external prediction accuracies up to 77%. A library of 240,000 potential CNT surface modifiers was further screened using these models and the least toxic organic modifiers were selected for experimental validation. Subsequent synthesis and testing of these surface-

modified CNTs confirmed the *in silico* predictions, demonstrating the utility of QNAR models for rational design of nanomaterials with enhanced properties.

In another study, a logistic regression-based QNAR model was developed²⁴⁷ to flag toxic outcomes; this model was trained on high-throughput toxicity screening data for BEAS2B cells exposed to nine metal oxide nanoparticles. The best-performing model had almost 100% classification accuracy and required only three nanoparticle descriptors: the period of the nanoparticle metal; the atomization energy of the metal oxide; and the nanoparticle size and volume fraction. Another study used RF classification to model cellular toxicity of metal oxide ENMs.²⁴⁸ The model was trained on data extracted from 216 publications, and used 14 ENMs attributes as descriptors. It demonstrated that cytotoxicity of ENMs was highly correlated with the administered dose, assay type, exposure time, and surface area of nanoparticles.²⁴⁸

Bayesian networks as models for predictive toxicology and for assessment of causal relationships

Models that predict toxicity of ENMs must account not only for the properties of the nanomaterials *per se*, but also for experimental conditions (e.g., assay types, exposure concentrations, exposure period, organism and more). It is important to quantify the relevance and significance of ENMs and experimental attributes driving toxicity while accounting for uncertainties in data, particularly that collected from multiple sources. Toxicity prediction models trained on these attribute combinations can sometimes identify causal relationships,²³⁹ which can be effectively achieved with the Bayesian Network (BN, also called a Bayesian Belief network, BBN) approach.²⁴⁹

BN models construct a network where the nodes are ENMs characteristics and the edges (links) represent conditional dependences of target outcomes on various attributes. This provides a visual representation of causal relationships.²⁵⁰ The model allows interpretation of “if/then” causal relationships where the parent (antecedent) and child (descendent) nodes are at the outgoing and incoming links in the BN structure, respectively. The set of model attributes and their conditional dependencies represents knowledge from the dataset(s) of attributes and toxicity outcomes in the form of probability distributions. BN models can identify, for example, the conditional dependence that would lead to a toxicity outcome within a specific range.

Previous studies have demonstrated the value of BNs for developing qualitative “toxicity/hazard” classification of ENMs based on using physicochemical and specialized descriptors.²⁵¹ BN models identified the most relevant parameters impacting specific ENMs hazards. Thus, regression and classification models were developed²⁵² for cause–effect relationships for hazard associated with exposure to TiO₂, SiO₂, Ag, CeO₂, ZnO NPs for different toxicity endpoints. A BN model predicted the hazard associated with exposure to metal and metal oxide NPs²⁵¹ for eight toxicity endpoints compiled from 32 published studies. Despite the existence of significant data gaps for some NPs the resulting BN model identified the most relevant NP properties for predicting toxicity outcomes.

Data variability and curation

As is true for traditional QSAR, inter- and intra-sample variability in QNAR is a big issue that can dramatically affect the predictivity of a model. Therefore, in order to study and/or model nanomaterials, the experimental variability for both inter- and intra-sample measurements needs to be taken into account whenever possible. For instance, the size distribution of a given sample of a specified nanomaterial can vary from one instrument to another. If a series of size distribution plots is used to model a set of nanomaterials, then the experimental variability of these measured profiles needs to be considered to better understand the stability, reliability, and robustness of the model. As with small molecule drugs and/or batches of biologics, replicate measurements are necessary to understand experimental variability. All, or a subset of compounds chosen to be representative, and their associated samples are characterized in triplicate. If one endpoint (e.g., particle diameter, zeta potential) is deemed unreliable, that endpoint should not be considered as a descriptor for those nanomaterials nor should it be considered as a target property for a model. Clearly, materials characteristics measured with low accuracy and reproducibility, will limit the predictivity of the QNAR models trained using them. Nanomaterials are particularly sensitive to the protocols used for sample preparations (e.g., dilution, sonication, solvent mixtures) leading to aggregation or even degradation. Experimental variability is a general issue that the QSAR modeling field is constantly dealing with. Strict data curation prior to model development is highly recommended,¹⁵ whereas external validation ensures the stability and robustness of the models over all modeling and external prediction sets.

Perspectives

Although QNAR modeling is still in its infancy, we anticipate it will grow significantly in the near future. This growth is dependent on:

- development of more effective and interpretable ENMs-specific descriptors
- further development of high-throughput synthesis and screening platforms for nanomaterials, leading to the expansion of publicly available data to train QNAR models
- development of more robust and predictive, consensus models based on individual QNAR models trained on diverse ENMs descriptors using advanced ML techniques including DL
- development of nanomaterials with desired properties and pre-computed bioprofiles generated by interdisciplinary research teams. The role of QNAR modeling in the context of such multidisciplinary efforts cannot be overestimated.

Biomaterials and regenerative medicine

Previous sections have covered major underlying concepts of cheminformatics such as chemical similarity, QSAR model building and validation, and domain of applicability. These methods have been progressively extended to areas beyond their traditional applications, for instance chemical genomics and (nano)materials science as discussed

above. Another emerging field is the use of QSAR methods to model control of cell phenotypes and understanding and predicting the biological response to materials. These are relatively recent, but rapidly expanding fields where the potential impact is very significant. Unlike bioinformatics,²⁵³ cell biology, and clinical medicine,²⁵⁴ there is a relative paucity of published examples of the application of QSAR or related ML-based methods to biomaterials, regenerative medicine, and stem cells studies. Polymers and other complex materials have been used in implantable or indwelling devices, as replacement or augmentation of natural bodily components, as scaffolds for cell culture, and as active biomaterials and drug delivery systems. Unfortunately, such materials are not as well defined as organic molecules. As discussed above in the sections on (nano)materials informatics, one of the biggest challenges in the field of biomaterials is generation of appropriate descriptors that capture relevant properties of these materials and can adequately represent their structure, often poorly understood and characterized.¹⁹ In this regard, rapid adoption of DL methods is providing useful models for this very important issue. The feature generation capabilities of DNN mean that simpler representations of complex materials become possible. We further anticipate that predictive material-QSAR models may be interrogated to identify the types of complex features that modulate relevant biological responses most strongly.

Although the use of arcane molecular descriptors has already resulted in good predictive models of the biological effects of materials, there is increasing impatience with their inability to be related back to underlying chemical features interpretable by chemists to improve performance. The dilemma between good predictions of properties for new materials, and interpretability of models (mechanistically or in terms of molecular interactions at a surface) has been reviewed recently by Fujita and Winkler.¹⁹ This nexus has led to a rise in the popularity of signature or fragment-based descriptors for modeling of materials interaction with biological systems. For example, signature descriptors have been used to model the adhesion of bacteria to polymers.²⁵⁵ New ML methods such as adversarial and encoder-decoder networks have begun tackling the ‘inverse QSAR’ problem, where trained model can be used to design or suggest new molecules for synthesis with improved activity.

A second important issue that distinguishes materials modeling from small molecules modelling is that in the former case interactions are more complex. Often materials interact with mixtures of proteins, membranes, cells, and modulate the responses of a myriad signal pathways, mechanosensors, *etc.* Consequently, ML methods are best suited to address such complexity and uncertainty, where the mechanisms of the cell-materials interactions are largely unknown. Notably, ML methods have been successfully used already for modeling soft biological materials such as blood vessels.²⁵⁶

To date, QSAR methodology has been applied in regenerative medicine and biomaterials modeling in three major groups of studies. First, sparse and non-sparse feature selection methods have been used to reduce the complexity of materials-biological systems interactions. For example, sparse feature selection methods were applied to investigate stem cell behavior (see Figure 8 for details). Similarly, an expectation maximization algorithm employing a sparse (Laplacian) prior⁴⁵ was used to identify the most relevant genes in

unbiased genome-wide expression studies. In one such study, mesenchymal stem cells (MSCs) were exposed to the components of a biomaterial (strontium bioglass, SrBG) with varying levels of strontium ions.²⁵⁷ These drive MSC differentiation down the osteogenic pathway to form bone tissue. After preliminary expression level and fold ratio filtering, the sparse feature selection method identified a handful of genes related to fatty and sterol biosynthesis - a previously unreported mechanism of bone growth modulation. Subsequent experimental validation of this mechanism by means of qPCR Raman spectroscopy and protein expression profiling led to important implications for the control osteoporosis and bone loss.

In another related investigation, unbiased sparse feature selection methods were applied to gene expression data.²⁵⁸ In this experiment, stem cells were forced to divide symmetrically or asymmetrically in response to several types of experimental conditions.²⁵⁸ Sparse feature selection methods were used to identify robust markers for symmetric cell division, which is a very important factor in stem cell proliferation and differentiation studies.²⁵⁸

ML methods have been increasingly applied to quantitative modeling of the responses of biological systems to interactions with materials.²⁵⁹ To date, most of these materials have been polymers, due to their tunable properties, ease of library generation and characterization, and generally understood biocompatibility. Early work was conducted by the Kohn group from Rutgers University who generated a library of 112 tyrosine-derived polyarylates and measured a range of their physical properties and biological responses.²⁶⁰ They used DRAGON descriptors²⁶¹ based on the monomeric units of the polymers in combination with such parameters as glass transition temperature (T_g) and air-water contact angle to generate quantitative and predictive models of fetal rat lung fibroblast (FRLF) metabolism and fibrinogen attachment on the polymer surfaces. Subsequently, research teams at the University of Nottingham, CSIRO, Monash University, and MIT generated polymer microarrays²⁶² and conducted high throughput screening to elucidate structure-property relationships in their interactions with cells.

The use of biomaterials as cell factories²⁶³ shows great promise, and the large generated stem cell attachment, proliferation, and differentiation datasets were modelled by ML methods. These could make robust and accurate predictions of stem cell behavior of materials not used to train the models. In one study, the attachment of embryoid bodies (a surrogate and stable cell system to mimic embryonic stem cells) to a polymer library was modelled using sparse feature selection and optimally regularized neural networks.²⁶⁴ These models relied on DRAGON descriptors and Bayesian regularized neural networks to quantify the attachment of embryoid bodies to the polyacrylate libraries. A more recent study modelled attachment, proliferation, and differentiation of human dental pulp stem cells to a polymer library.²⁶⁵ In this case study, the authors also investigated the ability of a 541 members of polyacrylate homopolymer and copolymer library to promote attachment, proliferation, and differentiation of stem cells.

Finally, advanced QSAR methods are being applied to the characterization of surfaces that interact with biological systems and to analyzes of complex high-content data such as cell imaging and phenotype recognition. Surface analysis methods such as Raman and Time-of-

Flight Secondary Ion Mass Spectrometry (ToF-SIMS) are invaluable experimental tools for characterizing the nature of surfaces interacting with biology. Surprisingly, there has been little application of statistical methods and ML to the corresponding spectroscopic data. ToF-SIMS in particular has proven to generate data that is very useful for QSPR material modeling.²⁶² Recent work has shown how self-organizing maps (SOMS) can provide superior clustering of complex mass peak data,²⁶⁶ probing into the intrinsic information content (Shannon entropy) of these surface analysis methods.²⁶⁷

As the field of biomaterials modeling is relatively nascent, there are many issues that need resolving before the full benefit of AI/ML-based QSAR methods can be realized. The most important of these issues is how to represent a high molecular weight complex material such as a cross-linked polymer hydrogel or polymer library with distributions of chain length, block sizes, degree of cross-linking, etc. Although surprisingly effective models can be generated using descriptors based on small fragments, additional materials features may be needed where these approximations fail. More recently methods have been developed that allow many types of nanoscale topographies to be imprinted onto materials surfaces. These modulate biological properties such as macrophage polarization, so efficient ways of generating descriptors for topographical features are required. Equally important is the need to generate models that can be interrogated to guide the synthesis of subsequent generations of materials with improved characteristics.¹⁷⁴ Biological data variability and reproducibility are also a constant struggle for high throughput materials-based experiments. Improving the reliability of these biological response data by careful statistical treatment of results and improved fabrication quality control is also important. However, as modeling of biomaterials coevolves with further development of the respective experimental research, one shall expect models to become more robust and impactful.

Clinical and health informatics

Just as advances in statistics, ML, and AI have influenced chemical research, experience accumulated in cheminformatics can be applied to clinical research. The growing linkage between QSAR modeling and clinical informatics was highlighted by the most recent 22nd EuroQSAR meeting in 2018 dedicated explicitly to “Translational and Health Informatics: Implications for Drug Discovery”.²⁶⁸ One example of such cross-fertilization between the fields is the development of robotic biomarkers of motor impairment of patients recovering from stroke.²⁶⁹

One of the greatest challenges in designing clinical trials is dealing with the subjectivity and variability introduced by human assessment of clinical endpoints. This problem is particularly acute in neurology, where outcomes may be highly variable (e.g., in cognition), susceptible to the state of the patient (e.g., fatigue, pain, anxiety, depression), the lack of a gold standard definition or diagnosis (e.g., neuropathy, dementia), are high dimensional (e.g., imaging or genomic markers), or are composite in nature (e.g., clinical instruments for assessing depression or quality of life).²⁷⁰ These factors make it difficult to demonstrate treatment benefits, requiring larger pools of subjects in clinical trials as well as properly structured electronic health record (EMR) archiving and retrieval capabilities.

Neurological disorders such as stroke suffer from clinical assessment limitations as established methods are often subjective: Scales such as the Fugl-Meyer (FM),²⁷¹ Motor Power (MP),²⁷² NIH Stroke (NIH),²⁷³ and Modified Rankin (MR),²⁷⁴ require properly trained personnel for evaluation, with results widely varying from rater to rater.²⁷⁵ While extensive training of raters and centralization of outcome assessments (whenever possible) can reduce variability, it does not completely eliminate it and comes with its own additional costs.²⁷⁶

One way to minimize this measurement variability issue is to replace human raters with robotic technology that can provide repeatable, reliable and speedy assessment of continuous measures of impairment and its change during recovery. Robotic devices are less sensitive to the skills and expertise of a human rater, can reduce inter- and intra-rater variability, can be used simultaneously for both assessment and rehabilitation, which can be done faster and more frequently, and can further be used in a home setting thus minimizing patient burden and inconvenience.²⁷⁶

The following study illustrates the use of QSAR -type approaches in clinical informatics. To test their utility in clinical trials, the four clinical scales mentioned above, were used in conjunction with a robotic assay to measure arm movement in 208 patients at 7, 14, 21, 30, and 90-day time-points after acute ischemic stroke. The data were collected at two clinical sites in the US and the UK. The study had two goals. The first was to establish whether the robotic measurements could predict the scores of human raters, and the second was to develop a more sensitive robotic biomarker that could reduce the sample size of the study without compromising the predictive value. The robots were low impedance and low friction interactive devices that measured speed, position, and force.²⁷⁷ The robotic assessment consisted of 35 macro- and micro-metrics derived from various directed, unassisted reaching, circle drawing, resistance to external forces, and shoulder strength measurements, applied to the affected and unaffected arms.²⁷⁸

The relationships between these 35 robotic variables and the four clinical scales were visualized (see Figure 9) using stochastic proximity embedding (SPE), a self-organizing nonlinear mapping algorithm that was originally invented to visualize very large combinatorial chemical libraries¹¹⁵ and subsequently adapted for various molecular modeling applications.²⁷⁹ Having established a degree of correlation, models were generated to assess whether the robotic metrics could predict the clinical scales with sufficient accuracy to serve as their surrogates. The model was trained using the data from degree of recovery from day 7 to day 90 after stroke, and all other intermediate measurements were used as test data. Specifically, 208 patients were divided into two complementary populations: those with complete data sets for days 7 and 90 (referred to as *completers*; N=87) and; those with missing data on days 7 or 90 (referred to as *non-completers*; N=121). The models, based on feed-forward NNs, were derived independently for each clinical scale. They were trained to predict the clinical scores of a given patient on a given day from the respective robotic metrics, using the *completer* population as a training set.

To minimize over-fitting, a feature selection algorithm based on artificial ant colonies, originally developed for QSAR applications, was used to identify the subset of robotic

metrics that had the highest predictive power.²⁸⁰ Once the relevant features were identified, ensemble models comprising 10 neural network predictors were constructed using the same network topology and training parameters but initialized with a different random number seed. The predictions of these models were averaged to produce an ensemble prediction. All models were cross-validated using the standard jackknife approach that divided the training data into 10 disjoint subsets containing 10% of the patterns each, systematically removing each subset from the training set, building a model with the remaining patterns, and predicting the clinical scores of the removed patterns using the optimized network parameters. The resulting predictions were compared to the original clinical scores to evaluate the overall agreement with the R^2_{CV} metrics. This process was repeated 10 times to obtain more robust cross-validation statistics. Finally, the best models identified by cross-validation were used to predict performance of the *non-completers*, who formed an independent test set. This protocol was virtually identical to the one used for QSAR applications.²⁸¹

The resulting models recapitulated the human scored clinical scales with a cross-validated R^2 of 0.73, 0.75, 0.63, and 0.60 for the FM, MP, NIH and MR scales, respectively. The models also showed lower but still useful predictive power for the external validation set (*non-completers*). The models had better prediction accuracy for the FM and MP scales that are more closely related to motor function than the NIH and MR metrics. Finally, the models were used to derive novel composite robotic endpoints with improved sensitivity (and effect size) compared to existing scales. To measure the effect size, Cohen's d parameter for paired observations was used, defined as the mean divided by the standard deviation of the day 7 to day 90 changes over all the *completers*. Since optimizing nonlinear composites is an ill-posed mathematical problem, a greedy forward-selection algorithm was employed to select up to 8 most relevant robotic features. Optimized robotic composites with as few as four features increased the effect size over a reference natural history trial²⁸² by as much as 107% for the training and 83% for the test set. This result is highly significant as an increase of 83% in effect size would result in a 70% reduction in the number of patients required to achieve the typical 80% statistical power in a clinical trial.

While the primary purpose of EMRs is to serve patient care, the second QSAR-inspired study illustrates how structured EMR information can be processed with unsupervised learning to improve patient phenotyping in Chronic Obstructive Pulmonary Disease (COPD).²⁸³ COPD, a heterogeneous disease characterized by persistent, non-reversible airflow limitation is the fourth leading cause of death in the United States (as of 2010). While "phenotype" is a co-emergent property of the genotype-environment interaction, COPD has been classically stratified in two phenotypes,²⁸⁴ the "blue bloater", which is rooted in chronic bronchitis (cyanosis due to hypoxemia), and the "pink puffer", which is rooted in emphysema (pink skin and hyperinflation), although up to seven COPD phenotypes have been proposed, based on "clinical relevance".²⁸⁵ Unsupervised learning was used to analyze EMR data from COPD patients, first to find out if common COPD patterns exist, which in turn could identify different COPD subtypes and lead to improved therapeutic management within each COPD subtype. A total of 3,144 patients aged 40 or older, admitted to the University of New Mexico Hospital, a 580-bed tertiary hospital with a COPD diagnosis (ICD9 codes: 490, 491, 492 or 496) between 1 January 2011 and 1 May

2014 were processed for this study. Data processed in this analysis included demographics, comorbidities, presence of atopy, obesity, number of admissions, prescriptions for inhalers (grouped as: i) short acting beta-agonist, ii) long-acting beta-agonist, iii) anticholinergics, iv) steroids and v) combinations), prescriptions for oral steroids, beta-blockers and statins, as well as weight loss and elevated plasma bicarbonate (used as surrogate biomarkers for disease severity). All variables, including age (40–65 years and >65 years) and number of admissions (one admission and \geq two admissions), were coded as binary for the study.

These data were clustered using the *sphere exclusion algorithm*,²⁸⁶ a disjoint similarity method that has been widely applied in cheminformatics. In the disjoint similarity method, a patient (object) can belong to only one cluster.²⁸⁷ When processing this multidimensional space that has as many dimensions as variables, dissimilarity can serve as the distance metric between patients. By definition, similarity is set to 0 if all the variables are different and is set to 1 if they are equal.²⁸⁷ As described elsewhere, in sphere exclusion the only user input is the similarity threshold: First, the similarity between all patients was computed. The algorithm then identified the patient with the most “neighbors” within a specified similarity cut-off, forming the first cluster. These patients were excluded from further iterations. The process was repeated until only patients without neighbors (i.e., singletons) were left. For this dataset, the optimal balance between the number of clusters and clustering overlap was found at similarity threshold 0.62. Using the sphere exclusion algorithm for clustering reduces the risk of bias since the method does not make *a priori* assumptions regarding numbers of clusters or similarity thresholds.

After leaving 189 patients (6%) as outliers, the following nine COPD clusters (phenotypes) were identified, with the number of patients given in brackets: 1: Depression–COPD (1748); 2: Malignancy–COPD (312); 3: Coronary artery disease–COPD (291); 4: Young age–low comorbidity–high readmission–COPD (152); 5: Advanced malignancy–COPD (144); 6: Cerebrovascular disease–COPD (120); 7: Atopy–COPD (81); 8: Diabetes mellitus – Chronic Kidney Disease – COPD (64) and 9: Advanced disease–COPD (43). The largest cluster is characterized by a large proportion of patients over age 65 and depression; two clusters (2 and 5) are associated with malignancy, although the first one has few readmissions whereas the second one has signs of advanced COPD and frequent readmissions. Cluster 3 is associated with heart disease (patients over age 65), whereas cluster 6 is associated with predominantly cerebrovascular disease and younger (under 65) patients. Cluster 4 (young patients, few comorbidities) has the highest number of prescriptions for bronchodilators; cluster 7 is also comprised of patients below age 65, but with asthma/atopy and higher numbers of readmissions; cluster 8 is associated with chronic kidney disease (CKD) and type 2 diabetes in patients aged 40–65, whereas cluster 9 has frequent readmissions, severe disease and high number of anticholinergic prescriptions. Our analysis revealed five previously unreported COPD phenotypes: two malignancy-COPD clusters (2 and 5), the COPD – CKD – diabetes cluster (8), the “advanced disease” cluster (9) and the high readmission phenotype (4). Each of these new clusters has practical implications, which may lead to better therapeutic outcomes.

To summarize, the above studies successfully adapted methods from computational chemistry and cheminformatics into in-depth analyses of health data. We anticipate that this

transfer of methods and experience will continue to fuel healthcare informatics research by introducing new and improved computational methodologies.

Outlook

The field of QSAR modeling based on simple approaches used to predict chemical reactivity was initially popularized by Corwin Hansch and his colleagues more than 55 years ago.¹ For many years, even decades, this field was focused on the prediction of physicochemical properties and biological activities using descriptors representing intrinsic properties of chemical structures. However, as the size and diversity of chemical datasets expanded, the QSAR modeling field has evolved to include larger and more diverse types of chemical descriptors and increasingly more complex statistical and machine learning techniques. We reflected on these trends earlier,² and foreshadowed the impact that these developments in the QSAR modeling community would have on many other areas of research. We projected that, with the continuing strong growth of publicly accessible data, this field will become essential for extracting knowledge from, and making predictions with, these massive data sets. We forecast that the field will continue to embrace even more powerful and complex machine learning methods. Furthermore, we expect that these modeling methods will continue to find rapid acceptance not only in chemistry but also in new fields beyond chemistry, where large data sets are readily available and modeling complex relationships between a set of independent variables and given properties of interest are important. The recent expansion of QSAR studies using DL approaches (as discussed in the section on Modern trends in QSAR) is an early harbinger of these expectations.

We have illustrated some of non-traditional applications in this review, demonstrating how QSAR-like approaches are beginning to yield exciting results in research areas as diverse as quantum mechanics, materials and nanomaterials science, biomaterials, regenerative medicine, and health care. Impressively, many of the roadblocks and technical issues in statistical data modelling employed in different domains of knowledge had already been addressed in the QSAR modeling literature. Examples include papers on the impact of the errors on QSAR analysis²⁸⁸ and the importance of data curation to achieve stable and reproducible models.⁹³ These considerations were under active discussion in the QSAR community before the reproducibility crisis brought to light by the NIH²⁸⁹ and biomedical scientific community at large.²⁹⁰ Similarly, rigorous model validation prior to prediction¹⁴ and the importance of rigor in modeling protocols²⁹¹ have been articulated in several seminal publications in QSAR field²⁹² and have already been adopted as regulatory requirements.⁸⁴ Extreme examples of the application of QSAR concepts beyond its traditional domain are provided by a study into factors influencing temporal crime patterns in Chicago²⁹³ that cites a well-known work on QSAR model validation²⁹² and a study on stock price predictions.²⁹⁴

We expect QSAR-like modeling techniques to continue to expand substantially even beyond the areas where it is starting to make an impact, which we discussed above. Scientists working in this field will continue to experiment with novel statistical, machine learning, and AI algorithms to accelerate the experimental discovery of novel compounds and materials with desired properties. The jury is still out on whether the newest DL approaches

will improve the prediction accuracy of QSAR models. However, we expect that the answer will emerge in the next few years, given the tremendous activity in this field.

As discussed above, stunning and potentially paradigm shifting developments are occurring in the use of machine learning approaches to massively accelerate quantum mechanical calculations, without sacrificing accuracy, and the use of QSAR methods for *de novo* compound design. Another fascinating and emerging direction is AI-driven chemical synthesis route prediction and its synergy with robotic synthesis, also discussed above. We anticipate a multitude of new and interesting algorithmic developments in the area of retro- and forward synthesis design, with software integrated with the robotic systems. We should soon see the emergence of fully autonomous, 'close loop' chemical and materials synthesis and optimization systems. In addition to these methodological developments, we foresee many new and impactful experimental methods arising that lead to novel, useful, and safe chemicals when QSAR modeling is applied to these data, and the increased application of ML methodologies in drug target selection, gene-phenotype evaluation and disease modeling. Finally, besides potentially exciting developments in traditional areas of application in chemical sciences, we further expect that the experience in model development, validation, and exploitation of QSAR models for knowledge discovery in chemical sciences will lead to progressive expansion of QSAR modeling principles and approaches in many other disciplines.

Conclusions

This contribution was conceived by a group of scientists who have dedicated significant portions of their professional careers to the development and use of quantitative methods in computational chemistry and molecular modeling. Following the previous highly cited comprehensive survey of QSAR modeling that was coauthored by many contributors to this paper and published in 2014,² we felt it was time to reflect on the new and exciting developments in QSAR modeling that have emerged in the last five years due to proliferation of large and diverse (Big Data) molecular bioactivity datasets and of burgeoning use of associated Big Data analytical methods such as DL. We also intended to share our observations and excitement concerning the prolific use of similar ML approaches in areas beyond chemical domain; the latter excitement and observations were in part influenced by the transition to other fields that some original cheminformaticians, including several coauthors of this paper, have made in their own research evolution and career development. Herein, we have summarized recent and developing trends in several areas of research where statistical data modeling has begun taking a prominent place and where experiences and generalizable approaches of QSAR modeling could catalyze new discoveries. We hope that this collective contribution will be useful for both specialists in data modeling and experimental researchers looking to expand their toolkits to include computational data analytical approaches.

Acknowledgements

This review combined a series of separately written, invited contributions from the various coauthors (some sections with multiple coauthors). Primary attributions for the various contributed sections are as follows: Introduction – E. Muratov and A. Tropsha; Chemical similarity – J. Bajorath; Modern trends in QSAR modeling – R. Sheridan;

QSAR in chemical safety assessment – I. Tetko; Multi-target profiling and polypharmacology – D. Filimonov and V. Porokov; QSAR-like approaches in chemical genomics – T. Oprea and A. Cherkasov; QSAR in synthetic organic chemistry – I. Baskin and A. Varnek; Closed-loop discovery and automation – A. Aspuru-Guzik; Machine learning approaches in quantum chemistry – O. Isayev and A. Roitberg; Materials informatics – S. Curtarolo; Nanomaterials informatics – Y. Cohen and D. Fourches; Biomaterials and regenerative medicine – D. Winkler; Clinical and health informatics – D. Agrafiotis and T. Oprea; Outlook and Conclusions – E. Muratov, A. Cherkasov, and A. Tropsha. Final editing was accomplished by A. Tropsha, who also takes primary responsibility for the final content. Mentioning of trade names or commercial products does not constitute endorsement or recommendation for use.

The authors acknowledge the seminal contributions of Corwin Hansch and Toshio Fujita in the initial development of the QSAR field and of Frank Burden in the development of sparse feature selection and Bayesian regularization of neural networks for QSAR. The authors acknowledge many fruitful discussions with members of their groups. TIO acknowledge NIH funding support (U24 CA224370, U24 TR002278, and U01CA239108). AT and EM acknowledge NIH funding support (U01CA207160). VP and DF would like to acknowledge the support of the Russian Program for Basic Research of State Academies of Sciences for 2013-2020.

Notes and references

- Hansch C, Maloney P, Fujita T and Muir R, *Nature*, 1962, 194, 178–180.
- Cherkasov A, Muratov EN, Fourches D, Varnek A, Baskin II, Cronin M, Dearden JC, Gramatica P, Martin YC, Todeschini R, Consonni V, Kuz'min VEVE, Cramer RD, Benigni R, Yang C, Rathman JF, Terflath L, Gasteiger J, Richard AM and Tropsha A, *J. Med. Chem.*, 2014, 57, 4977–5010. [PubMed: 24351051]
- Ban F, Dalal K, Li H, LeBlanc E, Rennie PS and Cherkasov A, *J. Chem. Inf. Model.*, 2017, 57, 1018–1028. [PubMed: 28441481]
- Alves VM, Muratov EN, Zakharov A, Muratov NN, Andrade CH and Tropsha A, *Food Chem. Toxicol.*, 2018, 112, 526–534. [PubMed: 28412406]
- Simón-Vidal L, García-Calvo O, Oteo U, Arrasate S, Lete E, Sotomayor N and González-Díaz H, *J. Chem. Inf. Model.*, 2018, 58, 1384–1396. [PubMed: 29898360]
- Sheridan R, Schafer W, Piras P, Zawatzky K, Sherer EC, Roussel C and Welch CJ, *J. Chromatogr. A*, 2016, 1467, 206–213. [PubMed: 27318509]
- Grzybowski BA, Szymku S, Gajewska EP, Molga K, Dittwald P, Wołos A and Klucznik T, *Chem.*, 2018, 4, 390–398.
- Capuzzi SJ, Sun W, Muratov EN, Martínez-Romero C, He S, Zhu W, Li H, Tawa G, Fisher EG, Xu M, Shinn P, Qiu X, García-Sastre A, Zheng W and Tropsha A, *J. Med. Chem.*, 2018, 61, 3582–3594. [PubMed: 29624387]
- Hong M, Chen X, Zhang R, Wang D, Shen S and Singh VP, *Ocean Sci.*, 2018, 14, 301–320.
- Ghosh D and Guha R, *Comput. Environ. Urban Syst.*, 2010, 34, 189–203.
- Muratov EN, Lewis M, Fourches D, Tropsha A and Cox WC, *Am. J. Pharm. Educ.*, 2017, 81, 46. [PubMed: 28496266]
- Hosseini R, Newlands N, Dean C, Takemura A, Hosseini R, Newlands NK, Dean CB and Takemura A, *Remote Sens.*, 2015, 7, 2752–2780.
- Oprea T, Olah M, Ostopovici L, Rad R and Mracec M, in *EuroQSAR 2002—Designing Drugs and Crop Protectants: Processes Problems and Solutions*, eds. Ford M, Livingstone D, Dearden J and Van de Waterbeemd HH, Blackwell Publishing, New York, 2003, pp. 314–315.
- Golbraikh A and Tropsha A, *J. Mol. Graph. Model.*, 2002, 20, 269–276. [PubMed: 11858635]
- Fourches D, Muratov E and Tropsha A, *J. Chem. Inf. Model.*, 2016, 56, 1243–1252. [PubMed: 27280890]
- Editorial, *Nature*, 2014, 515, 7–7.
- Tropsha A, *Mol. Inform.*, 2010, 29, 476–488. [PubMed: 27463326]
- Lowe D, In the pipeline, <https://blogs.sciencemag.org/pipeline/archives/2018/01/30/automated-chemistry-a-vision>, (accessed 19 August 2019).
- Fujita T and Winkler DA, *J. Chem. Inf. Model.*, 2016, 56, 269–274. [PubMed: 26754147]
- Peltason L, Iyer P and Bajorath J, *J. Chem. Inf. Model.*, 2010, 50, 1021–1033. [PubMed: 20443603]
- Peltason L and Bajorath J, *J. Med. Chem.*, DOI:10.1021/jm0705713.

22. Maggiora GM, *J. Chem. Inf. Model*, 2006, 46, 1535. [PubMed: 16859285]
23. Kosloff M and Kolodny R, *Proteins*, 2008, 71, 891–902. [PubMed: 18004789]
24. Bajorath J, Peltason L, Wawer M, Guha R, Lajiness MSMSMS and Van Drie JHJH, *Drug Discov. Today*, 2009, 14, 698–705. [PubMed: 19410012]
25. Schneider G, Neidhart W, Giller T and Schmid G, *Angew. Chemie Int. Ed*, 1999, 38, 2894–2896.
26. Lo Y-C, Rensi SE, Torng W and Altman RB, *Drug Discov. Today*, 2018, 23, 1538–1546. [PubMed: 29750902]
27. PARACELTUS TPAB, *Opera Omnia Medico-Chemico-Chirurgica, tribus voluminibus comprehensa., Sumptibus Joan. Antonii, & Samuelis De Tournes, Geneva, Editio 11., 1658.*
28. Lavecchia A, *Drug Discov. Today*, 2015, 20, 318–331. [PubMed: 25448759]
29. Gaulton A, Bellis LJ, Bento AP, Chambers J, Davies M, Hersey A, Light Y, McGlinchey S, Michalovich D, Al-Lazikani B and Overington JP, *Nucleic Acids Res.*, 2012, 40, D1100–7. [PubMed: 21948594]
30. 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 and Bryant SH, *Nucleic Acids Res.*, 2015, 44, D1202–13. [PubMed: 26400175]
31. Sheridan RP, *J. Chem. Inf. Model*, 2013, 53, 783–90. [PubMed: 23521722]
32. Muratov EN, Varlamova EV, Artemenko AG, Polishchuk PG and Kuz'min VE, *Mol. Inform*, 2012, 31, 202–221. [PubMed: 27477092]
33. Oprisiu I, Varlamova E, Muratov E, Artemenko A, Marcou G, Polishchuk P, Kuz'min V and Varnek A, *Mol. Inform*, 2012, 31, 491–502. [PubMed: 27477467]
34. Zakharov AV, Varlamova EV, Lagunin AA, Dmitriev AV, Muratov EN, Fourches D, Kuz'min VE, Poroikov VV, Tropsha A and Nicklaus MC, *Mol. Pharm*, 2016, 13, 545–556. [PubMed: 26669717]
35. Segall MD, *Curr. Drug Metab*, 2012, 18, 1292–1310.
36. Prado-Prado FJ, González-Díaz H, de la Vega OM, Ubeira FM and Chou K-C, *Bioorg. Med. Chem*, 2008, 16, 5871–80. [PubMed: 18485714]
37. Brown JB, Okuno Y, Marcou G, Varnek A and Horvath D, *J. Comput. Aided. Mol. Des*, 2014, 28, 597–618. [PubMed: 24771144]
38. Van Westen GJP, Wegner JK, Ijzerman AP, Van Vlijmen HWT and Bender A, *Medchemcomm*, 2011, 2, 16–30.
39. DREAM Challenges, IDG-DREAM Drug-Kinase Binding Prediction Challenge - Dream Challenges, <http://dreamchallenges.org/project/idg-dream-drug-kinase-binding-prediction-challenge/>, (accessed 1 January 2020).
40. Eriksson L, Jaworska J, Worth AP, Cronin MTD, McDowell RM and Gramatica P, *Environ. Health Perspect*, 2003, 111, 1361–1375. [PubMed: 12896860]
41. Golbraikh A, Muratov E, Fourches D and Tropsha A, *J. Chem. Inf. Model*, 2014, 54, 1–4. [PubMed: 24251851]
42. Cramer RD III, Patterson DE and Bunce JD, *J. Am. Chem. Soc*, 1988, 110, 5959–5967. [PubMed: 22148765]
43. Kuz'min VE, Muratov EN, Artemenko AG, Gorb L, Qasim M and Leszczynski J, *J. Comput. Aided. Mol. Des*, 2008, 22, 747–759. [PubMed: 18385948]
44. Polishchuk P, *J. Chem. Inf. Model*, 2017.
45. Burden FR and Winkler DA, *QSAR Comb. Sci*, 2009, 28, 645–653.
46. Artemenko A, Muratov E, Kuz'min V, Kovdienko N, Hromov A, Makarov V, Riabova O, Wutzler P and Schmidtke M, *J. Antimicrob. Chemother*, 2007, 60, 68–77. [PubMed: 17550890]
47. Polishchuk P, Kuz'min V, Artemenko A and Muratov E, *Mol. Inform*, 2013, 32, 843–853. [PubMed: 27480236]
48. Sheridan RP, *J. Chem. Inf. Model*, 2019, 59, 1324–1337. [PubMed: 30779563]
49. Hansch C, *Acc. Chem. Res*, 1993, 26, 147–153.
50. Mater AC and Coote ML, *J. Chem. Inf. Model*, 2019, 59, 2545–2559. [PubMed: 31194543]
51. Ma J, Sheridan RP, Liaw A, Dahl GE and Svetnik V, *J. Chem. Inf. Model*, 2015, 55, 263–274. [PubMed: 25635324]

52. MERCK, Kaggle Merck Molecular Activity Challenge, <https://www.kaggle.com/c/MerckActivity>, (accessed 19 August 2019).
53. Cover T and Hart P, *IEEE Trans. Inf. Theory*, 1967, 13, 21–27.
54. Wold S, Sjöström M and Eriksson L, *Chemom. Intell. Lab. Syst.*, 2001, 58, 109–130.
55. Geppert H, Horváth T, Gärtner T, Wrobel S and Bajorath J, *J. Chem. Inf. Model*, 2008, 48, 742–746. [PubMed: 18318473]
56. Burden FR and Winkler DA, *J. Chem. Inf. Model*, 2015, 55, 1529–1534. [PubMed: 26158341]
57. Breiman LEO, *Mach. Learn.*, 2001, 45, 5–32.
58. Dudley R, *J. Funct. Anal.*, 1967, 1, 290–330.
59. Svetnik V, Wang T, Tong C, Liaw A, Sheridan RP and Song Q, *J. Chem. Inf. Model*, 2005, 45, 786–799. [PubMed: 15921468]
60. Sheridan RP, *J. Chem. Inf. Model*, 2013, 53, 2837–2850. [PubMed: 24152204]
61. Sheridan RP, Wang WM, Liaw A, Ma J and Gifford EM, *J. Chem. Inf. Model*, 2016, 56, 2353–2360. [PubMed: 27958738]
62. Ke G, Meng Q, Finley T, Wang T, Chen W, Ma W, Ye Q and Liu T-Y, in *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS, Long Beach, 2017*, pp. 3149–3157.
63. Lenselink EB, ten Dijke N, Bongers B, Papadatos G, van Vlijmen HWT, Kowalczyk W, IJzerman AP and van Westen GJP, *J. Cheminform*, 2017, 9, 45. [PubMed: 29086168]
64. Winkler DA and Le TC, *Mol. Inform*, 2017, 36, 1600118.
65. Cybenko G, *Math. Control. Signals, Syst.*, 1989, 2, 303–314.
66. Golbraikh A, Fourches D, Sedykh A, Muratov E, Liepina I and Tropsha A, in *Practical Aspects of Computational Chemistry III*, eds. Leszczynski J and Shukla M, Springer, New York, Heidelberg, Dordrecht, London, 2014, pp. 187–230.
67. Hochreiter S, Klambauer G and Rarey M, *J. Chem. Inf. Model*, 2018, 58, 1723–1724. [PubMed: 30109927]
68. Ramsundar B, Liu B, Wu Z, Verras A, Tudor M, Sheridan RP and Pande V, *J. Chem. Inf. Model*, DOI:10.1021/acs.jcim.7b00146.
69. Varnek A, Gaudin C, Marcou G, Baskin I, Pandey AK and Tetko IV, *J. Chem. Inf. Model*, 2009, 49, 133–144. [PubMed: 19125628]
70. Xu Y, Ma J, Liaw A, Sheridan RP and Svetnik V, *J. Chem. Inf. Model*, 2017, 57, 2490–2504. [PubMed: 28872869]
71. Coley CW, Barzilay R, Green WH, Jaakkola TS and Jensen KF, *J. Chem. Inf. Model*, 2017, 57, 1757–1772. [PubMed: 28696688]
72. Faber FA, Hutchison L, Huang B, Gilmer J, Schoenholz SS, Dahl GE, Vinyals O, Kearnes S, Riley PF and von Lilienfeld OA, *J. Chem. Theory Comput*, 2017, 13, 5255–5264. [PubMed: 28926232]
73. Merk D, Friedrich L, Grisoni F and Schneider G, *Mol. Inform*, 2018, 37, 1700153.
74. Dacrema MF, Cremonesi P and Jannach D, in *Proceedings of the 13th ACM Conference on Recommender Systems - RecSys '19*, ACM Press, New York, 2019, pp. 101–109.
75. Bach S, Binder A, Montavon G, Klauschen F, Müller K-R and Samek W, *PLoS One*, 2015, 10, e0130140. [PubMed: 26161953]
76. Baskin II, Winkler D and Tetko IV, *Expert Opin. Drug Discov*, 2016, 11, 785–795. [PubMed: 27295548]
77. Garcia Denegri ME, Bustillo S, Gay CC, Van De Velde A, Gomez G, Echeverría S, Gauna Pereira MDC, Maruñak S, Nuñez S, Bogado F, Sanchez M, Teibler GP, Fusco L and Leiva LCA, *Curr. Top. Med. Chem.*, DOI:10.2174/1568026619666190725094851.
78. Myatt GJ, Ahlberg E, Akahori Y, Allen D, Amberg A, Anger LT, Aptula A, Auerbach S, Beilke L, Bellion P, Benigni R, Bercu J, Booth ED, Bower D, Brigo A, Burden N, Cammerer Z, Cronin MTD, Cross KP, Custer L, Dettwiler M, Dobo K, Ford KA, Fortin MC, Gad-McDonald SE, Gellatly N, Gervais V, Glover KP, Glowienke S, Van Gompel J, Gutsell S, Hardy B, Harvey JS, Hillegass J, Honma M, Hsieh J-H, Hsu C-W, Hughes K, Johnson C, Jolly R, Jones D, Kemper R, Kenyon MO, Kim MT, Kruhlak NL, Kulkarni SA, Kümmerer K, Leavitt P, Majer B, Masten S, Miller S, Moser J, Mumtaz M, Muster W, Neilson L, Oprea TI, Patlewicz G, Paulino A, Lo Piparo

- E, Powley M, Quigley DP, Reddy MV, Richarz A-N, Ruiz P, Schilter B, Serafimova R, Simpson W, Stavitskaya L, Stidl R, Suarez-Rodriguez D, Szabo DT, Teasdale A, Trejo-Martin A, Valentin J-P, Vuorinen A, Wall BA, Watts P, White AT, Wichard J, Witt KL, Woolley A, Woolley D, Zwickl C and Hasselgren C, *Regul. Toxicol. Pharmacol.*, 2018, 96, 1–17. [PubMed: 29678766]
79. Ankley GT, Bennett RS, Erickson RJ, Hoff DJ, Hornung MW, Johnson RD, Mount DR, Nichols JW, Russom CL, Schmieder PK, Serrano J, Tietge JE and Villeneuve DL, *Environ. Toxicol. Chem.*, 2010, 29, 730–41. [PubMed: 20821501]
80. Pittman ME, Edwards SW, Ives C and Mortensen HM, *Toxicol. Appl. Pharmacol.*, 2018, 343, 71–83. [PubMed: 29454060]
81. Rybacka A, Rudén C, Tetko IV and Andersson PL, *Chemosphere*, 2015, 139, 372–378. [PubMed: 26210185]
82. Wittwehr C, Aladjov H, Ankley G, Byrne HJ, de Knecht J, Heinzle E, Klambauer G, Landesmann B, Luijten M, MacKay C, Maxwell G, (Bette) Meek ME, Paini A, Perkins E, Sobanski T, Villeneuve D, Waters KM and Whelan M, *Toxicol. Sci.*, 2017, 155, 326–336. [PubMed: 27994170]
83. US EPA, Tox21, <http://www.epa.gov/ncct/Tox21/>, (accessed 20 August 2019).
84. Organisation for Economic Co-operation and Development and OECD, OECD principles for the validation, for regulatory purposes, of (Quantitative) Structure-Activity Relationship models, <http://europa.eu.int/comm/environment/chemicals/reach.htm>, (accessed 17 February 2017).
85. Huang R, Xia M, Nguyen D-T, Zhao T, Sakamuru S, Zhao J, Shahane SA, Rossoshek A and Simeonov A, *Front. Environ. Sci.*, 2016, 3, 85.
86. Mayr A, Klambauer G, Unterthiner T and Hochreiter S, *Front. Environ. Sci.*, 2016, 3, 80.
87. Tetko IV, in *Methods in molecular biology*, Humana Press, Clifton, 2008, vol. 458, pp. 180–197.
88. Wu Y and Wang G, *Int. J. Mol. Sci.*, 2018, 19, 2358.
89. Mansouri K, Abdelaziz A, Rybacka A, Roncaglioni A, Tropsha A, Varnek A, Zakharov A, Worth A, Richard AM, Grulke CM, Trisciuzzi D, Fourches D, Horvath D, Benfenati E, Muratov E, Wedebe EB, Grisoni F, Mangiatori GF, Incisivo GM, Hong H, Ng HW, Tetko IV, Balabin I, Kancherla J, Shen J, Burton J, Nicklaus M, Cassotti M, Nikolov NG, Nicolotti O, Andersson PL, Zang Q, Politi R, Beger RD, Todeschini R, Huang R, Farag S, Rosenberg SA, Slavov S, Hu X and Judson RS, *Environ. Health Perspect.*, 2016, 124, 1023–1033. [PubMed: 26908244]
90. Wang Z, Gerstein M and Snyder M, *Nat. Rev. Genet.*, 2009, 10, 57–63. [PubMed: 19015660]
91. Liu R, Yu X and Wallqvist A, *J. Cheminform.*, 2015, 7, 4. [PubMed: 25717346]
92. Novotarskyi S, Abdelaziz A, Sushko Y, Körner R, Vogt J and Tetko IV, *Chem. Res. Toxicol.*, 2016, 29, 768–775. [PubMed: 27120770]
93. Fourches D, Muratov E and Tropsha A, *J. Chem. Inf. Model.*, 2010, 50, 1189–204. [PubMed: 20572635]
94. Jamei M, *Curr. Pharmacol. Reports*, 2016, 2, 161–169.
95. Wetmore BA, Wambaugh JF, Ferguson SS, Sochaski MA, Rotroff DM, Freeman K, Clewell HJ, Dix DJ, Andersen ME, Houck KA, Allen B, Judson RS, Singh R, Kavlock RJ, Richard AM and Thomas RS, *Toxicol. Sci.*, 2012, 125, 157–174. [PubMed: 21948869]
96. Yamane J, Aburatani S, Imanishi S, Akanuma H, Nagano R, Kato T, Sone H, Ohsako S and Fujibuchi W, *Nucleic Acids Res.*, 2016, 44, 5515–5528. [PubMed: 27207879]
97. Abdelaziz A, Sushko Y, Novotarskyi S, Körner R, Brandmaier S and Tetko IV, *Comb. Chem. High Throughput Screen.*, 2015, 18, 420–38. [PubMed: 25747436]
98. Sosnin S, Karlov D, Tetko IV and Fedorov MV, *J. Chem. Inf. Model.*, 2019, 59, 1062–1072. [PubMed: 30589269]
99. Alves VM, Muratov EN, Capuzzi SJ, Politi R, Low Y, Braga RC, Zakharov AV, Sedykh A, Mokshyna E, Farag S, Andrade CH, Kuz'min VE, Fourches D and Tropsha A, *Green Chem.*, 2016, 18, 4348–4360. [PubMed: 28503093]
100. Low YS, Alves VM, Fourches D, Sedykh A, Andrade CH, Muratov EN, Rusyn I and Tropsha A, *J. Chem. Inf. Model.*, 2018, 58, 2203–2213. [PubMed: 30376324]
101. Montavon G, Samek W and Müller K-R, *Digit. Signal Process.*, 2018, 73, 1–15.
102. Roth BL, Sheffler DJ and Kroeze WK, *Nat. Rev. Drug Discov.*, 2004, 3, 353–359. [PubMed: 15060530]

103. Lagunin A, Filimonov D and Poroikov V, *Curr. Pharm. Des.*, 2010, 16, 1703–1717. [PubMed: 20222853]
104. Ivanov SM, Lagunin AA and Poroikov VV, *Drug Discov. Today*, 2016, 21, 58–71. [PubMed: 26272036]
105. Tarasova OA, Urusova AF, Filimonov DA, Nicklaus MC, Zakharov AV and Poroikov VV, *J. Chem. Inf. Model*, 2015, 55, 1388–1399. [PubMed: 26046311]
106. Lagunin A, Stepanchikova A, Filimonov D and Poroikov V, *Bioinformatics*, 2000, 16, 747–8. [PubMed: 11099264]
107. Filimonov DA, Poroikov VV, Karaicheva EI, Kazarian RK, Budunova AP, Mikhaïlovskii EM, Rudnitskikh AV, Goncharenko LV and Burov IV, *Eksp. Klin. Farmakol.*, 1995, 58, 56–62.
108. Pogodin PV, Lagunin AA, Rudik AV, Filimonov DA, Druzhilovskiy DS, Nicklaus MC and Poroikov VV, *Front. Chem.*, 2018, 6, 133. [PubMed: 29755970]
109. González-Díaz H, Arrasate S, Gómez-SanJuan A, Sotomayor N, Lete E, Besada-Porto L and Ruso JM, *Curr. Top. Med. Chem.*, 2013, 13, 1713–41. [PubMed: 23889050]
110. Martin YC, Kofron JL and Traphagen LM, *J. Med. Chem.*, 2002, 45, 4350–4358. [PubMed: 12213076]
111. Sheridan RP and Kearsley SK, *Drug Discov. Today*, 2002, 7, 903–11. [PubMed: 12546933]
112. Keiser MJ, Roth BL, Armbruster BN, Ernsberger P, Irwin JJ and Shoichet BK, *Nat. Biotechnol.*, 2007, 25, 197–206. [PubMed: 17287757]
113. Luo M, Wang XS, Roth BL, Golbraikh A and Tropsha A, *J. Chem. Inf. Model*, 2014, 54, 634–47. [PubMed: 24410373]
114. Luo H, Chen J, Shi L, Mikailov M, Zhu H, Wang K, He L and Yang L, *Nucleic Acids Res.*, 2011, 39, W492–W498. [PubMed: 21558322]
115. Agrafiotis DK, Lobanov VS and Salemme FR, *Nat. Rev. Drug Discov.*, 2002, 1, 337–346. [PubMed: 12120409]
116. Gupta-Ostermann D and Bajorath J, *F1000Research*, 2014, 3, 113. [PubMed: 25383183]
117. March-Vila E, Pinzi L, Sturm N, Tinivella A, Engkvist O, Chen H and Rastelli G, *Front. Pharmacol.*, 2017, 8, 298. [PubMed: 28588497]
118. Lapinsh M, Prusis P, Gutcaits A, Lundstedt T and Wikberg JE, *Biochim. Biophys. Acta*, 2001, 1525, 180–90. [PubMed: 11342268]
119. Lapins M, Worachartcheewan A, Spjuth O, Georgiev V, Prachayasittikul V, Nantasenamat C and Wikberg JES, *PLoS One*, 2013, 8, e66566. [PubMed: 23799117]
120. Paricharak S, Cortés-Ciriano I, IJzerman AP, Malliavin TE and Bender A, *J. Cheminform.*, 2015, 7, 15. [PubMed: 25926892]
121. Orchard S, Al-Lazikani B, Bryant S, Clark D, Calder E, Dix I, Engkvist O, Forster M, Gaulton A, Gilson M, Glen R, Grigorov M, Hammond-Kosack K, Harland L, Hopkins A, Larminie C, Lynch N, Mann RK, Murray-Rust P, Lo Piparo E, Southan C, Steinbeck C, Wishart D, Hermjakob H, Overington J and Thornton J, *Nat. Rev. Drug Discov.*, 2011, 10, 661–9. [PubMed: 21878981]
122. Fourches D, Muratov E and Tropsha A, *Nat. Chem. Biol.*, 2015, 11, 535. [PubMed: 26196763]
123. Reymond J-L, *Acc. Chem. Res.*, 2015, 48, 722–730. [PubMed: 25687211]
124. Oprea TI, Bologa CG, Edwards BS, Prossnitz ER and Sklar LA, *J. Biomol. Screen.*, 2005, 10, 419–26. [PubMed: 16093551]
125. The Gene Ontology Consortium, *Nucleic Acids Res.*, 2017, 45, D331–D338. [PubMed: 27899567]
126. Hsing M, Byler K and Cherkasov A, *BMC Syst. Biol.*, 2008, 2, 80. [PubMed: 18796161]
127. Sedykh A, Zhu H, Tang H, Zhang L, Richard A, Rusyn I and Tropsha A, *Environ. Health Perspect.*, 2011, 119, 364–370. [PubMed: 20980217]
128. Bologa CG, Ursu O, Halip L, Curp n R and Oprea TI, *Rev. Roum. du Chim.*, 2015, 60, 219–226.
129. Woo G, Fernandez M, Hsing M, Lack NA, Cavga AD and Cherkasov A, *Bioinformatics*, 2019, 10.1093/bioinformatics/btz645.
130. Himmelstein DS and Baranzini SE, *PLOS Comput. Biol.*, 2015, 11, e1004259. [PubMed: 26158728]

131. The UniProt Consortium, *Nucleic Acids Res.*, 2017, 45, D158–D169. [PubMed: 27899622]
132. Kanehisa M, Furumichi M, Tanabe M, Sato Y and Morishima K, *Nucleic Acids Res.*, 2017, 45, D353–D361. [PubMed: 27899662]
133. Chen T and Guestrin C, in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '16*, ACM Press, New York, New York, USA, 2016, pp. 785–794.
134. Agarwal P and Searls DB, *Nat. Rev. Drug Discov.*, 2009, 8, 865–878. [PubMed: 19876041]
135. Nguyen D-T, Mathias S, Bologna C, Brunak S, Fernandez N, Gaulton A, Hersey A, Holmes J, Jensen LJ, Karlsson A, Liu G, Ma'ayan A, Mandava G, Mani S, Mehta S, Overington J, Patel J, Rouillard AD, Schürer S, Sheils T, Simeonov A, Sklar LA, Southall N, Ursu O, Vidovic D, Waller A, Yang J, Jadhav A, Oprea TI and Guha R, *Nucleic Acids Res.*, 2017, 45, D995–D1002. [PubMed: 27903890]
136. Oprea TI, Bologna CG, Brunak S, Campbell A, Gan GN, Gaulton A, Gomez SM, Guha R, Hersey A, Holmes J, Jadhav A, Jensen LJ, Johnson GL, Karlson A, Leach AR, Ma'ayan A, Malovannaya A, Mani S, Mathias SL, McManus MT, Meehan TF, von Mering C, Muthas D, Nguyen D-T, Overington JP, Papadatos G, Qin J, Reich C, Roth BL, Schürer SC, Simeonov A, Sklar LA, Southall N, Tomita S, Tudose I, Ursu O, Vidovic D, Waller A, Westergaard D, Yang JJ and Zahoránszky-Köhalmi G, *Nat. Rev. Drug Discov.*, 2018, 17, 317–332. [PubMed: 29472638]
137. Elsevier. 2018. “Reaxys Fact Sheet.”
138. Lowe DM, Glen R and Murray-Rust P, University of Cambridge, 2012.
139. Lin AI, Madzhidov TI, Klimchuk O, Nugmanov RI, Antipin IS and Varnek A, *J. Chem. Inf. Model.*, 2016, 56, 2140–2148. [PubMed: 27783508]
140. Liu B, Ramsundar B, Kawthekar P, Shi J, Gomes J, Luu Nguyen Q, Ho S, Sloane J, Wender P and Pande V, *ACS Cent. Sci.*, 2017, 3, 1103–1113. [PubMed: 29104927]
141. Polishchuk P, Madzhidov T, Gimadiev T, Bodrov A, Nugmanov R and Varnek A, *J. Comput. Aided. Mol. Des.*, 2017, 31, 829–839. [PubMed: 28752345]
142. Patel H, Bodkin MJ, Chen B and Gillet VJ, *J. Chem. Inf. Model.*, 2009, 49, 1163–1184. [PubMed: 19382767]
143. Hoonakker F, Lachiche N, Varnek A and Wagner A, *Int. J. Artif. Intell. Tools.*, 2010, 20, 253–270.
144. Varnek A, Fourches D, Hoonakker F and Solov'ev VP, *J. Comput. Aided. Mol. Des.*, 2005, 19, 693–703. [PubMed: 16292611]
145. Kowalik M, Gothard CM, Drews AM, Gothard NA, Weckiewicz A, Fuller PE, Grzybowski BA and Bishop KJM, *Angew. Chemie Int. Ed.*, 2012, 51, 7928–7932.
146. Chen L and Gasteiger J, *Angew. Chemie Int. Ed. English*, 1996, 35, 763–765.
147. Corey EJ, *Chem. Soc. Rev.*, 1988, 17, 111–133.
148. Segler MHS and Waller MP, *Chem. – A Eur. J.*, 2017, 23, 5966–5971.
149. Wei JN, Duvenaud D and Aspuru-Guzik A, *ACS Cent. Sci.*, 2016, 2, 725–732. [PubMed: 27800555]
150. Kayala MA, Azencott C-A, Chen JH and Baldi P, *J. Chem. Inf. Model.*, 2011, 51, 2209–2222. [PubMed: 21819139]
151. Szymku S, Gajewska EP, Klucznik T, Molga K, Dittwald P, Startek M, Bajczyk M and Grzybowski BA, *Angew. Chemie Int. Ed.*, 2016, 55, 5904–5937.
152. Coley CW, Green WH and Jensen KF, *Acc. Chem. Res.*, 2018, 51, 1281–1289. [PubMed: 29715002]
153. Karpov P, Godin G and Tetko IV, Springer, Cham, 2019, pp. 817–830.
154. Hartenfeller M, Zettl H, Walter M, Rupp M, Reisen F, Proschak E, Weggen S, Stark H and Schneider G, *PLoS Comput Biol.*, 2012, 8, e1002380. [PubMed: 22359493]
155. Méndez-Lucio O and Medina-Franco JL, *Drug Discov. Today*, 2017, 22, 120–126. [PubMed: 27575998]
156. Ertl P and Schuffenhauer A, *J. Cheminform.*, 2009, 1, 8. [PubMed: 20298526]
157. Coley CW, Rogers L, Green WH and Jensen KF, *J. Chem. Inf. Model.*, 2018, 58, 252–261. [PubMed: 29309147]

158. Taft RW, *J. Am. Chem. Soc.*, 1952, 74, 3120–3128.
159. Nugmanov RI, Madzhidov TI, Khaliullina GR, Baskin II, Antipin IS and Varnek AA, *J. Struct. Chem.*, 2014, 55, 1026–1032.
160. Glavatskikh M, Madzhidov T, Horvath D, Nugmanov R, Gimadiev T, Malakhova D, Marcou G and Varnek A, *Mol. Inform.*, 2019, 38, 1800077.
161. Gimadiev TR, Madzhidov TI, Nugmanov RI, Baskin II, Antipin IS and Varnek A, *J. Comput. Aided. Mol. Des.*, 2018, 32, 401–414. [PubMed: 29380104]
162. Marcou G, Aires de Sousa J, Latino DARS, de Luca A, Horvath D, Rietsch V and Varnek A, *J. Chem. Inf. Model.*, 2015, 55, 239–250. [PubMed: 25588070]
163. Gao H, Struble TJ, Coley CW, Wang Y, Green WH and Jensen KF, *ACS Cent. Sci.*, 2018, 4, 1465–1476. [PubMed: 30555898]
164. Hoonakker F, Lachiche N, Varnek A and Wagner A, in *Trends in Applied Intelligent Systems, Pt II, Proceedings*, Springer, Berlin, Heidelberg, 2010, vol. 6097, pp. 318–326.
165. Neri D and Lerner RA, *Annu. Rev. Biochem.*, 2018, 87, 479–502. [PubMed: 29328784]
166. Nikolaev P, Hooper D, Webber F, Rao R, Decker K, Krein M, Poleski J, Barto R and Maruyama B, *npj Comput. Mater.*, 2016, 2, 16031.
167. Saikin SK, Kreisbeck C, Sheberla D, Becker JS and A. A-G, *Expert Opin. Drug Discov.*, 2019, 14, 1–4. [PubMed: 30488727]
168. Häse F, Roch LM, Kreisbeck C and Aspuru-Guzik A, *ACS Cent. Sci.*, 2018, 4, 1134–1145. [PubMed: 30276246]
169. Häse F, Roch LM and Aspuru-Guzik A, *Chem. Sci.*, 2018, 9, 7642–7655. [PubMed: 30393525]
170. Roch LM, Häse F, Kreisbeck C, Tamayo-Mendoza T, Yunker LPE, Hein JE and Aspuru-Guzik A, *Sci. Robot.*, 2018, 3, eaat5559. [PubMed: 33141686]
171. Tabor DP, Roch LM, Saikin SK, Kreisbeck C, Sheberla D, Montoya JH, Dwaraknath S, Aykol M, Ortiz C, Tribukait H, Amador-Bedolla C, Brabec CJ, Maruyama B, Persson KA and Aspuru-Guzik A, *Nat. Rev. Mater.*, 2018, 3, 5–20.
172. Gómez-Bombarelli R, Aguilera-Iparraguirre J, Hirzel TD, Duvenaud D, Maclaurin D, Blood-Forsythe MA, Chae HS, Einzinger M, Ha D-G, Wu T, Markopoulos G, Jeon S, Kang H, Miyazaki H, Numata M, Kim S, Huang W, Hong SI, Baldo M, Adams RP and Aspuru-Guzik A, *Nat. Mater.*, 2016, 15, 1120–1127. [PubMed: 27500805]
173. Häse F, Roch LM and Aspuru-Guzik A, *Trends Chem.*, 2019, 1, 282–291.
174. Le TC and Winkler DA, *Chem. Rev.*, 2016, 116, 6107–6132. [PubMed: 27171499]
175. Butler KT, Davies DW, Cartwright H, Isayev O and Walsh A, *Nature*, 2018, 559, 547–555. [PubMed: 30046072]
176. Rupp M, Tkatchenko A, Müller K-R and von Lilienfeld OA, *Phys. Rev. Lett.*, 2012, 108, 058301. [PubMed: 22400967]
177. Hansen K, Biegler F, Ramakrishnan R, Pronobis W, von Lilienfeld OA, Müller K-R and Tkatchenko A, *J. Phys. Chem. Lett.*, 2015, 6, 2326–2331. [PubMed: 26113956]
178. Yao K, Herr JE, Brown SN and Parkhill J, *J. Phys. Chem. Lett.*, 2017, 8, 2689–2694. [PubMed: 28573865]
179. Huang B and Von Lilienfeld OA, *J. Chem. Phys.*, 2016, 145, 161102. [PubMed: 27802646]
180. Yao K, Herr JE and Parkhill J, *J. Chem. Phys.*, 2017, 146, 014106. [PubMed: 28063436]
181. Keenan AB, Jenkins SL, Jagodnik KM, Koplev S, He E, Torre D, Wang Z, Dohlan AB, Silverstein MC, Lachmann A, Kuleshov MV, Ma'ayan A, Stathias V, Terryn R, Cooper D, Forlin M, Koleti A, Vidovic D, Chung C, Schurer SC, Vasiliauskas J, Pilarczyk M, Shamsaei B, Fazal M, Ren Y, Niu W, Clark NA, White S, Mahi N, Zhang L, Kouril M, Reichard JF, Sivaganesan S, Medvedovic M, Meller J, Koch RJ, Birtwistle MR, Iyengar R, Sobie EA, Azeloglu EU, Kaye J, Osterloh J, Haston K, Kalra J, Finkbiener S, Li J, Milani P, Adam M, Escalante-Chong R, Sachs K, Lenail A, Ramamoorthy D, Fraenkel E, Daigle G, Hussain U, Coye A, Rothstein J, Sareen D, Ornelas L, Banuelos M, Mandefro B, Ho R, Svendsen CN, Lim RG, Stocksdale J, Casale MS, Thompson TG, Wu J, Thompson LM, Dardov V, Venkatraman V, Matlock A, Van Eyk JE, Jaffe JD, Papanastasiou M, Subramanian A, Golub TR, Erickson SD, Fallahi-Sichani M, Hafner M, Gray NS, Lin JR, Mills CE, Muhlich JL, Niepel M, Shamu CE, Williams EH, Wrobel D, Sorger

- PK, Heiser LM, Gray JW, Korkola JE, Mills GB, LaBarge M, Feiler HS, Dane MA, Bucher E, Nederlof M, Sudar D, Gross S, Kilburn DF, Smith R, Devlin K, Margolis R, Derr L, Lee A and Pillai A, *Cell Syst*, 2018, 6, 13–24. [PubMed: 29199020]
182. Schütt KT, Saucedo HE, Kindermans PJ, Tkatchenko A and Müller KR, *J. Chem. Phys.*, 2018, 148, 241722. [PubMed: 29960322]
183. Bartók AP, Kondor R and Csányi G, *Phys. Rev. B - Condens. Matter Mater. Phys.*, 2013, 87, 1–16.
184. De S, Bartók AP, Csányi G and Ceriotti M, *Phys. Chem. Chem. Phys.*, 2016, 18, 13754. [PubMed: 27101873]
185. Gastegger M, Schwiedrzik L, Bittermann M, Berzsenyi F and Marquetand P, *J. Chem. Phys.*, 2018, 148, 241709. [PubMed: 29960372]
186. Schütt KT, Arbabzadah F, Chmiela S, Müller KR and Tkatchenko A, *Nat. Commun.*, 2017, 8, 13890. [PubMed: 28067221]
187. Behler J and Parrinello M, *Phys. Rev. Lett.*, 2007, 98, 146401. [PubMed: 17501293]
188. Gastegger M, Behler J and Marquetand P, *Chem. Sci.*, 2017, 8, 6924–6935. [PubMed: 29147518]
189. Smith JS, Isayev O and Roitberg AE, *Chem. Sci.*, 2017, 8, 3192–3203. [PubMed: 28507695]
190. Fink T, Bruggesser H and Reymond JL, *Angew. Chemie - Int. Ed.*, 2005, 44, 1504–1508.
191. Smith JS, Nebgen B, Lubbers N, Isayev O and Roitberg AE, *J. Chem. Phys.*, 2018, 148, 241733. [PubMed: 29960353]
192. Schütt KT, Kindermans P-J, Saucedo HE, Chmiela S, Tkatchenko A and Müller K-R, 2017, 1, 1–10.
193. Ramakrishnan R, Dral PO, Rupp M and von Lilienfeld OA, *Sci. data*, 2014, 1, 140022. [PubMed: 25977779]
194. Blum LC and Reymond J-L, *J. Am. Chem. Soc.*, 2009, 131, 8732–8733. [PubMed: 19505099]
195. Brauer B, Kesharwani MK, Kozuch S and Martin JML, *Phys. Chem. Chem. Phys.*, 2016, 18, 20905–20925. [PubMed: 26950084]
196. Pulvermüller F, *Nat. Rev. Neurosci.*, DOI:10.1038/2201358a0.
197. Caruana R, *Mach. Learn.*, 1997, 28, 41–75.
198. Curtarolo S, Hart GLW, Nardelli MB, Mingo N, Sanvito S and Levy O, *Nat. Mater.*, 2013, 12, 191–201. [PubMed: 23422720]
199. Maddox J, *Nature*, 1988, 335, 201–201.
200. Hautier G, Fischer C, Ehlacher V, Jain A and Ceder G, *Inorg. Chem.*, 2011, 50, 656–663. [PubMed: 21142147]
201. Perim E, Lee D, Liu Y, Toher C, Gong P, Li Y, Simmons WN, Levy O, Vlassak JJ, Schroers J and Curtarolo S, *Nat. Commun.*, 2016, 7, 12315. [PubMed: 27480126]
202. Ward L, O’Keefe SC, Stevick J, Jelbert GR, Aykol M and Wolverton C, *Acta Mater.*, 2018, 159, 102–111.
203. Ghiringhelli LM, Vybiral J, Levchenko SV, Draxl C and Scheffler M, *Phys. Rev. Lett.*, 2015, 114, 105503. [PubMed: 25815947]
204. Isayev O, Fourches D, Muratov EN, Osés C, Rasch K, Tropsha A and Curtarolo S, *Chem. Mater.*, 2015, 27, 735–743.
205. Stanev V, Osés C, Kusne AG, Rodriguez E, Paglione J, Curtarolo S and Takeuchi I, *npj Comput. Mater.*, 2018, 4, 29.
206. Walsh A, *Nat. Chem.*, 2015, 7, 274–275. [PubMed: 25803462]
207. Curtarolo S, Morgan D, Persson K, Rodgers J and Ceder G, *Phys. Rev. Lett.*, 2003, 91, 135503. [PubMed: 14525315]
208. Rupp M, Tkatchenko A, Müller K-R and von Lilienfeld OA, *Phys. Rev. Lett.*, 2012, 108, 058301. [PubMed: 22400967]
209. Bartók AP, De S, Poelking C, Bernstein N, Kermode JR, Csányi G and Ceriotti M, *Sci. Adv.*, 2017, 3, e1701816. [PubMed: 29242828]
210. Pilia G, Mannodi-Kanakkithodi A, Uberuaga BP, Ramprasad R, Gubernatis JE and Lookman T, *Sci. Rep.*, 2016, 6, 19375. [PubMed: 26783247]

211. de Jong M, Chen W, Notestine R, Persson K, Ceder G, Jain A, Asta M and Gamst A, *Sci. Rep.*, 2016, 6, 34256. [PubMed: 27694824]
212. Madsen GKH, *J. Am. Chem. Soc.*, 2006, 128, 12140–12146. [PubMed: 16967963]
213. Legrain F, Carrete J, van Roekeghem A, Curtarolo S and Mingo N, *Chem. Mater.*, 2017, 29, 6220–6227.
214. Carrete J, Mingo N, Wang S and Curtarolo S, *Adv. Funct. Mater.*, 2014, 24, 7427–7432.
215. Sanvito S, Oses C, Xue J, Tiwari A, Zic M, Archer T, Tozman P, Venkatesan M, Coey M and Curtarolo S, *Sci. Adv.*, 2017, 3, e1602241. [PubMed: 28439545]
216. Yong J, Jiang Y, Usanmaz D, Curtarolo S, Zhang X, Li L, Pan X, Shin J, Takeuchi I and Greene RL, *Appl. Phys. Lett.*, 2014, 105, 222403.
217. Oses C, Gossett E, Hicks D, Rose F, Mehl MJ, Perim E, Takeuchi I, Sanvito S, Scheffler M, Lederer Y, Levy O, Toher C and Curtarolo S, *J. Chem. Inf. Model.*, 2018, 58, 2477–2490. [PubMed: 30188699]
218. Körner W, Krugel G, Urban DF and Elsässer C, *Scr. Mater.*, 2018, 154, 295–299.
219. Möller JJ, Körner W, Krugel G, Urban DF and Elsässer C, *Acta Mater.*, 2018, 153, 53–61.
220. Duvenaud DK, Maclaurin D, Iparraguirre J, Bombarell R, Hirzel T, Aspuru-Guzik A and Adams RP, in *Advances in Neural Information Processing Systems 28*, ed. Cortes C, Curran Associates, Inc, New Yourk, 2015, pp. 2224–2232.
221. Widom M, *J. Mater. Res.*, 2018, 33, 2881–2898.
222. Lederer Y, Toher C, Vecchio KS and Curtarolo S, *Acta Mater.*, 2018, 159, 364–383.
223. Emery AA, Saal JE, Kirklin S, Hegde VI and Wolverton C, *Chem. Mater.*, 2016, 28, 5621–5634.
224. Cao Y, Fatemi V, Fang S, Watanabe K, Taniguchi T, Kaxiras E and Jarillo-Herrero P, *Nature*, 2018, 556, 43–50. [PubMed: 29512651]
225. Tawfik SA, Isayev O, Stampfl C, Shapter J, Winkler DA and Ford MJ, *Adv. Theory Simulations*, 2019, 2, 1800128.
226. Jain A, Hautier G, Ong SP and Persson K, *J. Mater. Res.*, 2016, 31, 977–994.
227. Muratov EN, Artemenko AG, Varlamova EV, Polischuk PG, Lozitsky VP, Fedchuk AS, Lozitska RL, Gridina TL, Koroleva LS, Sil'nikov VN, Galabov AS, a Makarov V, Riabova OB, Wutzler P, Schmidtk M and Kuz'min VE, *Future Med. Chem.*, 2010, 2, 1205–26. [PubMed: 21426164]
228. Rose F, Toher C, Gossett E, Oses C, Nardelli MB, Fornari M and Curtarolo S, *Comput. Mater. Sci.*, 2017, 137, 362–370.
229. Fouches D, Pu D and Tropsha A, *Comb. Chem. High Throughput Screen*, 2011, 14, 217–225. [PubMed: 21275889]
230. Fouches D, Barnes J, Day NC, Bradley P, Reed JZ and Tropsha A, *Chem. Res. Toxicol.*, 2010, 23, 171–83. [PubMed: 20014752]
231. Gajewicz A, *Nanoscale*, 2017, 9, 8435–8448. [PubMed: 28604902]
232. Haase H and Klaessig A, *EU US Roadmap Nanoinformatics 2030*, 2018.
233. Turku Centre for Biotechnology, NanoMiner, <https://bioscience.fi/>, (accessed 1 September 2019).
234. [Nanoinfo.org](https://nanoinfo.org), NanoDatabank, <https://nanoinfo.org/nanodatabank/>, (accessed 2 September 2019).
235. Sansone S-A, Rocca-Serra P, Field D, Maguire E, Taylor C, Hofmann O, Fang H, Neumann S, Tong W, Amaral-Zettler L, Begley K, Booth T, Bougueleret L, Burns G, Chapman B, Clark T, Coleman L-A, Copeland J, Das S, de Daruvar A, de Matos P, Dix I, Edmunds S, Evelo CT, Forster MJ, Gaudet P, Gilbert J, Goble C, Griffin JL, Jacob D, Kleinjans J, Harland L, Haug K, Hermjakob H, Sui SJH, Laederach A, Liang S, Marshall S, McGrath A, Merrill E, Reilly D, Roux M, Shamu CE, Shang CA, Steinbeck C, Trefethen A, Williams-Jones B, Wolstencroft K, Xenarios I and Hide W, *Nat. Genet.*, 2012, 44, 121. [PubMed: 22281772]
236. Marchese Robinson RL, Lynch I, Peijnenburg W, Rumble J, Klaessig F, Marquardt C, Rauscher H, Puzyn T, Purian R, Aberg C, Karcher S, Vriens H, Hoet P, Hoover MD, Hendren CO and Harper SL, *Nanoscale*, 2016, 8, 9919–9943. [PubMed: 27143028]
237. Ehrenberg MS, Friedman AE, Finkelstein JN, Oberdörster G and McGrath JL, *Biomaterials*, 2009, 30, 603–610. [PubMed: 19012960]

238. Shaw SY, Westly EC, Pittet MJ, Subramanian A, Schreiber SL and Weissleder R, *Proc. Natl. Acad. Sci. U. S. A.*, 2008, 105, 7387–7392. [PubMed: 18492802]
239. Oh E, Liu R, Nel A, Gemill KB, Bilal M, Cohen Y and Medintz IL, *Nat. Nanotechnol.*, 2016, 11, 479–86. [PubMed: 26925827]
240. Puzyn T, Rasulev B, Gajewicz A, Hu X, Dasari TP, Michalkova A, Hwang H-M, Toropov A, Leszczynska D and Leszczynski J, *Nat. Nanotechnol.*, 2011, 6, 175–8. [PubMed: 21317892]
241. Fourches D, Pu D, Tassa C, Weissleder R, Shaw SY, Mumper RJ and Tropsha A, *ACS Nano*, 2010, 4, 5703–5712. [PubMed: 20857979]
242. Ojha PK, Kar S, Roy K and Leszczynski J, *Nanotoxicology*, 2018, 1–21. [PubMed: 29251527]
243. Kovalishyn V, Abramenko N, Kopernyk I, Charochkina L, Metelytsia L, Tetko IV, Peijnenburg W and Kustov L, *Food Chem. Toxicol.*, 2018, 112, 507–517. [PubMed: 28802948]
244. Alves VM, Hwang D, Muratov E, Sokolsky-Papkov M, Varlamova E, Vinod N, Lim C, Andrade CH, Tropsha A and Kabanov A, *Sci. Adv.*, 2019, 5, eaav9784. [PubMed: 31249867]
245. Trinh TX, Choi J-S, Jeon H, Byun H-G, Yoon T-H and Kim J, *Chem. Res. Toxicol.*, 2018, 31, 183–190. [PubMed: 29439565]
246. Fourches D, Pu D, Li L, Zhou H, Mu Q, Su G, Yan B and Tropsha A, *Nanotoxicology*, 2016, 10, 374–83. [PubMed: 26525350]
247. Liu R, Rallo R, George S, Ji Z, Nair S, Nel AE and Cohen Y, *Small*, 2011, 7, 1118–1126. [PubMed: 21456088]
248. Ha MK, Trinh TX, Choi JS, Maulina D, Byun HG and Yoon TH, *Sci. Rep.*, 2018, 8, 3141. [PubMed: 29453389]
249. Money ES, Barton LE, Dawson J, Reckhow KH and Wiesner MR, *Sci. Total Environ.*, 2014, 473–474, 685–691.
250. Neapolitan RE, *Mol. Biol.*, 2003, 6, 674.
251. Marvin HJP, Bouzembrak Y, Janssen EM, van der Zande M, Murphy F, Sheehan B, Mullins M and Bouwmeester H, *Nanotoxicology*, 2017, 11, 123–133. [PubMed: 28044458]
252. Murphy F, Sheehan B, Mullins M, Bouwmeester H, Marvin HJP, Bouzembrak Y, Costa AL, Das R, Stone V and Tofail SAM, *Nanoscale Res. Lett.*, 2016, 11, 503. [PubMed: 27848238]
253. Cheng C and Worzel WP, in *Genetic Programming Theory and Practice XII*, eds. Riolo R, Worzel WP and Kotanchek M, 2014, pp. 1–15.
254. Molina E, Uriarte E, Santana L, Matos M and Borges F, *Curr. Bioinform.*, 2013, 8, 438–451.
255. Mikulskis P, Alexander MR and Winkler DA, *Adv. Intell. Syst.*, 2019, 1900045.
256. Cilla M, Pérez-Rey I, Martínez MA, Peña E and Martínez J, *Int. j. numer. method. biomed. eng.*, DOI:10.1002/cnm.3121.
257. Autefage H, Gentleman E, Littmann E, Hedegaard MAB, Von Erlach T, O'Donnell M, Burden FR, Winkler DA and Stevens MM, *Proc. Natl. Acad. Sci.*, 2015, 112, 4280–4285. [PubMed: 25831522]
258. Huh YH, Noh M, Burden FR, Chen JC, Winkler DA and Sherley JL, *Stem Cell Res.*, 2015, 14, 144–154. [PubMed: 25636161]
259. Hook AL, Anderson DG, Langer R, Williams P, Davies MC and Alexander MR, *Biomaterials*, 2010.
260. Smith JR, Kholodovych V, Knight D, Welsh WJ and Kohn J, *QSAR Comb. Sci.*, 2005, 24, 99–113.
261. Todeschini R and Consonni V, *Handbook of Molecular Descriptors*, Wiley-VCH Verlag GmbH, Weinheim, Germany, 2000, vol. 11.
262. Hook AL, Chang CY, Yang J, Luckett J, Cockayne A, Atkinson S, Mei Y, Bayston R, Irvine DJ, Langer R, Anderson DG, Williams P, Davies MC and Alexander MR, *Nat. Biotechnol.*, DOI:10.1038/nbt.2316.
263. Celiz AD, Smith JGW, Langer R, Anderson DG, Winkler DA, Barrett DA, Davies MC, Young LE, Denning C and Alexander MR, *Nat. Mater.*, 2014, 13, 570–579. [PubMed: 24845996]
264. Epa VC, Yang J, Mei Y, Hook AL, Langer R, Anderson DG, Davies MC, Alexander MR and Winkler DA, *J. Mater. Chem.*, 2012, 22, 20902–20906. [PubMed: 24092955]

265. Rasi Ghaemi S, Delalat B, Gronthos S, Alexander MR, Winkler DA, Hook AL and Voelcker NH, *ACS Appl. Mater. Interfaces*, 2018, 10, acsami.8b12473.
266. Madiona RMT, Bamford SE, Winkler DA, Muir BW and Pigram PJ, *Anal. Chem*, 2018, 90, 12475–12484. [PubMed: 30260219]
267. Madiona RMT, Welch NG, Russell SB, Winkler DA, Scoble JA, Muir BW and Pigram PJ, *Surf. Interface Anal.*, 2018, 50, 713–728.
268. 22nd EuroQSAR — Discngine - Enhancing Life Science Research, <https://www.discngine.com/events/2018/9/16/22nd-euroqsar>, (accessed 1 January 2020).
269. Krebs HI, Krams M, Agrafiotis DK, DiBernardo A, Chavez JC, Littman GS, Yang E, Byttebier G, Dipietro L, Rykman A, McArthur K, Hajjar K, Lees KR and Volpe BT, *Stroke*, 2014, 45, 200–204. [PubMed: 24335224]
270. Evans SR, *J. Exp. Stroke Transl. Med*, 2010, 3, 19–27. [PubMed: 21533012]
271. Fugl-Meyer AR, Jääskö L, Leyman I, Olsson S and Steglind S, *Scand. J. Rehabil. Med.*, 1975, 7, 13–31. [PubMed: 1135616]
272. Gregson J, Leathley MJ, Moore AP, Smith TL, Sharma AK and Watkins CL, *Age Ageing*, 2000, 29, 223–228. [PubMed: 10855904]
273. NIH, NIH Stroke Scale, <https://www.stroke.nih.gov/resources/scale.htm>, (accessed 29 August 2019).
274. Rankin J, *Scott. Med. J.*, 1957, 2, 200–215. [PubMed: 13432835]
275. Krebs HI, Volpe BT, Ferraro M, Fasoli S, Palazzolo J, Rohrer B, Edelstein L and Hogan N, *Top. Stroke Rehabil.*, 2002, 8, 54–70. [PubMed: 14523730]
276. Lo AC, Guarino PD, Richards LG, Haselkorn JK, Wittenberg GF, Federman DG, Ringer RJ, Wagner TH, Krebs HI, Volpe BT, Bever CT, Bravata DM, Duncan PW, Corn BH, Maffucci AD, Nadeau SE, Conroy SS, Powell JM, Huang GD and Peduzzi P, *N. Engl. J. Med.*, 2010, 362, 1772–1783. [PubMed: 20400552]
277. Hogan N and Krebs HI, in *Progress in brain research*, 2011, vol. 192, pp. 59–68. [PubMed: 21763518]
278. Flash T and Hogan N, *J. Neurosci.*, 1985, 5, 1688–703. [PubMed: 4020415]
279. Zhu F and Agrafiotis DK, *J. Comput. Chem.*, 2007, 28, 1234–1239. [PubMed: 17299775]
280. and SI and Agrafiotis D, , DOI:10.1021/CI000336S.
281. Agrafiotis DK, Cedeno W and Lobanov VS, *J. Chem. Inf. Comput. Sci.*, 2002, 42, 903–911. [PubMed: 12132892]
282. Kerr DM, Fulton RL, Lees KR and VISTA Collaborators, *Stroke*, 2012, 43, 1401–1403. [PubMed: 22308254]
283. Vazquez Guillamet R, Ursu O, Iwamoto G, Moseley PL and Oprea T, *Health Informatics J.*, 2018, 24, 394–409. [PubMed: 27856785]
284. Burrows B, Fletcher CM, Heard BE, Jones NL and Wootliff JS, *Lancet (London, England)*, 1966, 1, 830–5.
285. Mirza S and Benzo R, *Mayo Clin. Proc.*, 2017, 92, 1104–1112. [PubMed: 28688465]
286. Taylor R, *J. Chem. Inf. Model.*, 1995, 35, 59–67.
287. MacCuish J, Nicolaou C and MacCuish NE, *J. Chem. Inf. Comput. Sci.*, 2001, 41, 134–146. [PubMed: 11206366]
288. Young D, Martin D, Venkatapathy R, Harten P, Martin T, Venkatapathy R and Harten P, *QSAR Comb. Sci.*, 2008, 27, 1337–1345.
289. Collins FS and Tabak LA, *Nature*, 2014, 505, 612–3. [PubMed: 24482835]
290. Baker M, *Nature*, 2016, 533, 452–454. [PubMed: 27225100]
291. Dearden JC, Cronin MTD and Kaiser KLE, *SAR QSAR Environ. Res.*, 2009, 20, 241–266. [PubMed: 19544191]
292. Tropsha A, Gramatica P and Gombar VK, *QSAR Comb. Sci.*, 2003, 22, 69–77.
293. Towers S, Chen S, Malik A and Ebert D, *PLoS One*, 2018, 13, e0205151. [PubMed: 30356321]
294. Sheelapriya G and Murugesan R, *Spanish J. Financ. Account. / Rev. Española Financ. y Contab.*, 2017, 46, 189–211.

295. PICLIN N, PINTORE M, LANZA CM, SCACCO A, GUCCIONE S, GIURATO L and CHRÉTIEN JR, *J. Sens. Stud*, 2008, 23, 558–569.
296. Schut AGT, Stephens DJ, Stovold RGH, Adams M and Craig RL, *Crop Pasture Sci.*, 2009, 60, 60–70.
297. Xiao M and Obbard JP, *GCB Bioenergy*, 2010, 2, 346–352.
298. Alavi AH, Gandomi AH, Modaresnezhad M and Mousavi M, *J. Earthq. Eng.*, 2011, 15, 511–536.
299. Antelio M, Esteves MGP, Schneider D and de Souza JM, in 2012 IEEE International Conference on Systems, Man, and Cybernetics (SMC), IEEE, 2012, pp. 931–936.
300. Mousavi SM, Mostafavi ES and Hosseinpour F, *Comput. Ind. Eng.*, 2014, 74, 120–128.
301. Cao Y, Jiang Y, Gao H, Chen H, Fang X, Mu H and Tao F, *Comput. Electron. Agric.*, 2014, 106, 49–55.
302. El Haddad J, Canioni L and Bousquet B, *Spectrochim. Acta Part B At. Spectrosc.*, 2014, 101, 171–182.
303. Dearden JC, Cronin MT and Kaiser KL, *SAR QSAR Environ. Res.*, 2009, 20, 241–266. [PubMed: 19544191]
304. Ponomarenko J, Dizhbite T, Lauberts M, Viksna A, Dobele G, Bikovens O and Telysheva G, *BioResources*, 2014, 9, 2051–2068.
305. Sattar AMA, *J. Hydroinformatics*, 2014, 16, 550–571.
306. Elhakeem M and Sattar AMA, *Earth Surf. Process. Landforms*, 2015, 40, 1216–1226.
307. Tajeri S, Sadrossadat E and Bazaz JB, *Int. J. Rock Mech. Min. Sci.*, 2015, 80, 107–117.
308. Mundava C, Schut AGT, Helmholz P, Stovold R, Donald G and Lamb DW, *Rangel. J.*, 2015, 37, 157.
309. Heitzig S, Weinebeck A and Murrenhoff H, *SAE Int. J. Fuels Lubr.*, 2015, 8, 2015-01–9075.
310. Pan T-T, Sun D-W, Cheng J-H and Pu H, *Compr. Rev. Food Sci. Food Saf.*, 2016, 15, 529–541. [PubMed: 33401821]
311. Malaj E, Guénard G, Schäfer RB and von der Ohe PC, *Ecol. Appl.*, 2016, 26, 1249–59. [PubMed: 27509762]
312. Nikolaides A, Miess S, Auvera I, Müller R, Klosterkötter J and Ruhrmann S, *Eur. Arch. Psychiatry Clin. Neurosci.*, 2016, 266, 649–661. [PubMed: 27305925]
313. Polanski J, Kucia U, Duszkiwicz R, Kurczyk A, Magdziarz T and Gasteiger J, *Sci. Rep.*, 2016, 6, 28521. [PubMed: 27334348]
314. Tavana M, Fallahpour A, Di Caprio D and Santos-Arteaga FJ, *Expert Syst. Appl.*, 2016, 61, 129–144.
315. Ising HK, Ruhrmann S, Burger NAFM, Rietdijk J, Dragt S, Klaassen RMC, van den Berg DPG, Nieman DH, Boonstra N, Linszen DH, Wunderink L, Smit F, Veling W and van der Gaag M, *Psychol. Med.*, 2016, 46, 1839–1851. [PubMed: 26979398]
316. Sattar AMA, Gharabaghi B and McBean EA, *Water Resour. Manag.*, 2016, 30, 1635–1651.
317. Alavi AH, Hasni H, Zaabar I and Lajnef N, *Arch. Civ. Mech. Eng.*, 2017, 17, 326–335.
318. Mousavi SM, Mostafavi ES and Jiao P, *Energy Convers. Manag.*, 2017, 153, 671–682.
319. Hamze-Ziabari SM and Yasavoli A, *J. Adv. Concr. Technol.*, 2017, 15, 644–661.
320. SHAHRARA N, ÇELIK T and GANDOMI AH, *J. Civ. Eng. Manag.*, 2017, 23, 85–95.
321. Atieh M, Taylor G, Sattar AMA and Gharabaghi B, *J. Hydrol.*, 2017, 545, 383–394.
322. Tesfahunegn GB and Wortmann CS, *Nutr. Cycl. Agroecosystems*, 2017, 109, 269–289.
323. Cabrero JM and Yurrita M, *Eng. Struct.*, 2018, 171, 895–910.
324. Hou E, Wang J and Chen W, *Geocarto Int.*, 2018, 33, 754–769.
325. Kovdienko N, Polishchuk P, Muratov E, Artemenko A, Kuz'min V, Gorb L, Hill F and Leszczynski J, *Mol. Inform.*, 2010, 29, 394–406. [PubMed: 27463195]
326. Zhang X, Li X, Li L, Zhang S and Qin Q, *J. Arid Land*, DOI:10.1007/s40333-018-0110-2.
327. Najafzadeh M, Rezaie-Balf M and Tafarjnoruz A, *Int. J. River Basin Manag.*, 2018, 16, 505–512.
328. Haidl T, Rosen M, Schultze-Lutter F, Nieman D, Eggers S, Heinimaa M, Juckel G, Heinz A, Morrison A, Linszen D, Salokangas R, Klosterkötter J, Birchwood M, Patterson P, Ruhrmann S

and European Prediction of Psychosis Study (EPOS) Group, *Schizophr. Res.*, 2018, 199, 346–352. [PubMed: 29661524]

329. Glawe H, Sanna A, Gross EKV and Marques MAL, *New J. Phys.*, 2016, 18, 093011.
330. Isayev O, Oses C, Toher C, Gossett E, Curtarolo S and Tropsha A, *Nat. Commun.*, 2017, 8, 15679. [PubMed: 28580961]
331. Ouyang R, Curtarolo S, Ahmetcik E, Scheffler M and Ghiringhelli LM, *Phys. Rev. Mater.*, 2018, 2, 083802.

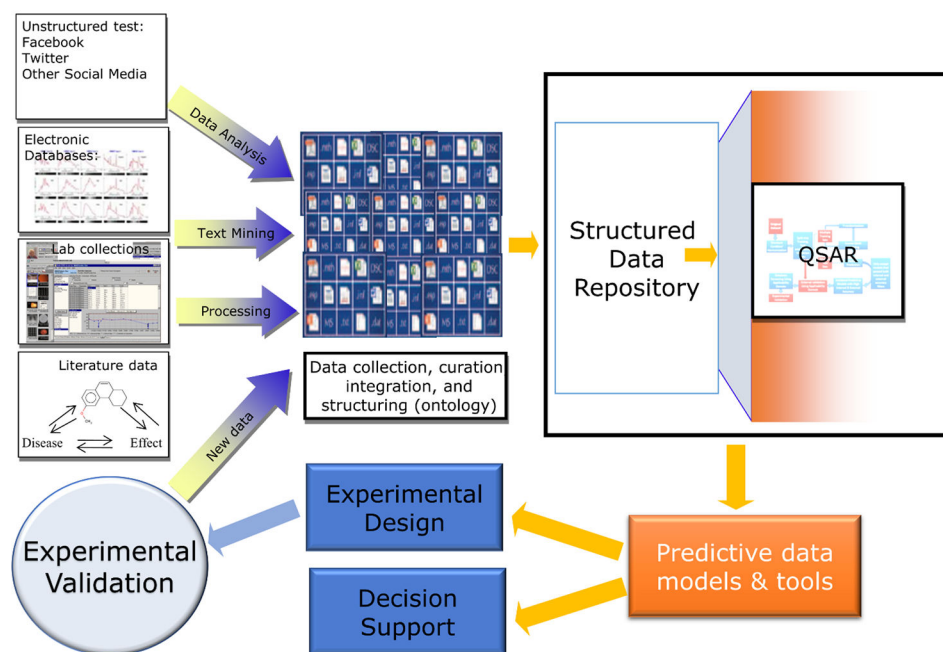


Figure 1. Data cycle associated with QSAR modeling projects.

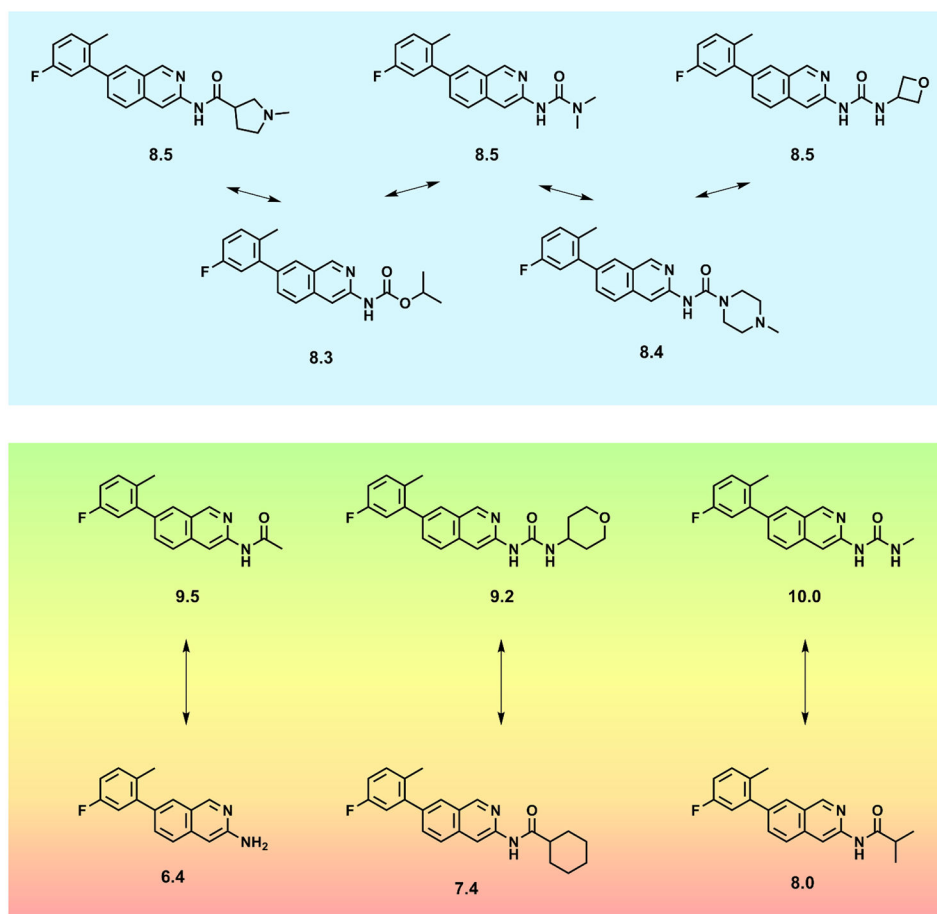


Figure 2. Different SAR patterns. Shown are inhibitors of tyrosine kinase ABL forming different SARs. For each compound the logarithmic potency (pKi) value is reported. At the top, SAR continuity is observed where gradually changes in compound structure (traced by horizontal arrows) are accompanied by moderate potency alterations. By contrast, the inhibitors at the bottom display SAR discontinuity. Here, small structural modifications lead to large changes in potency. Vertical arrows indicate the formation of pairwise activity cliffs.

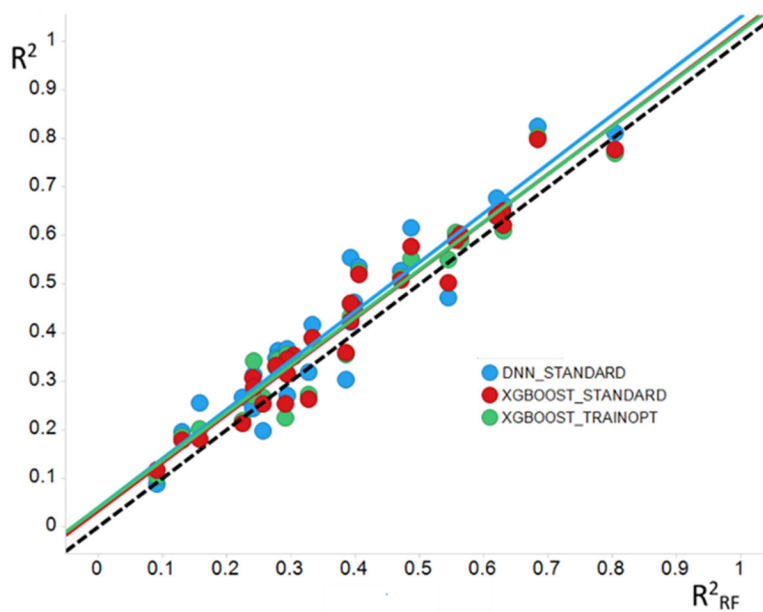


Figure 3. Comparison of the Pearson R² values for models generated using DNN (blue) or XGBoost (red and green) and random forest methods.

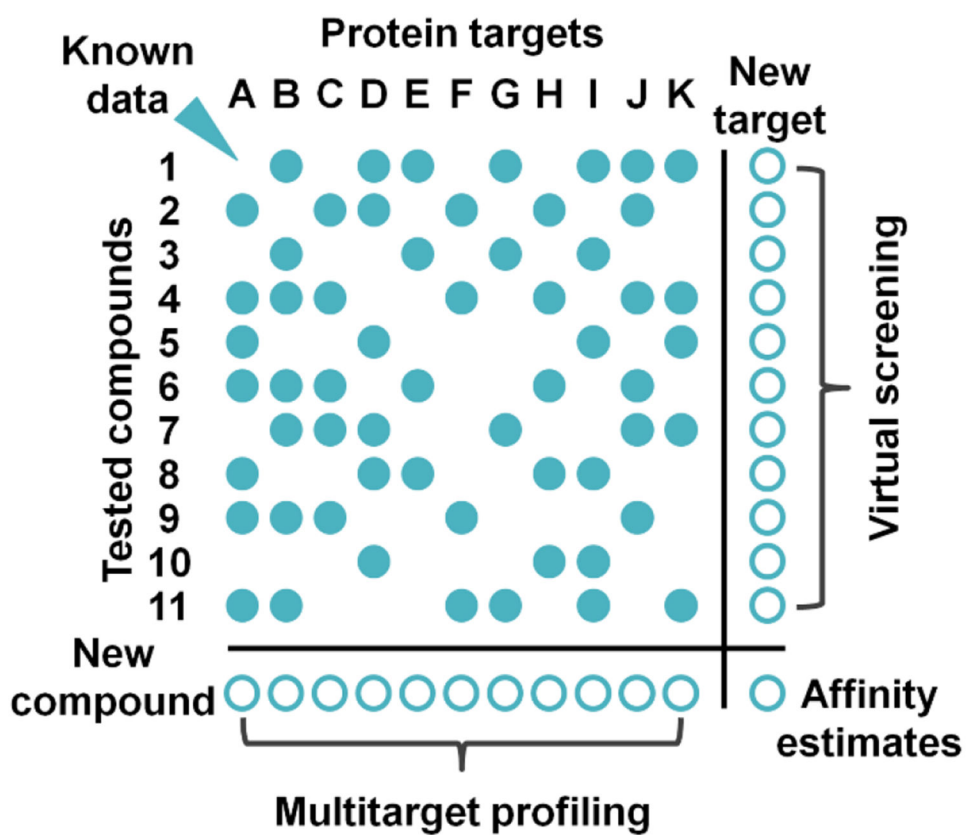


Figure 4. Proteochemometrics approach enables accurate affinity estimates for novel ligand-target pairs.

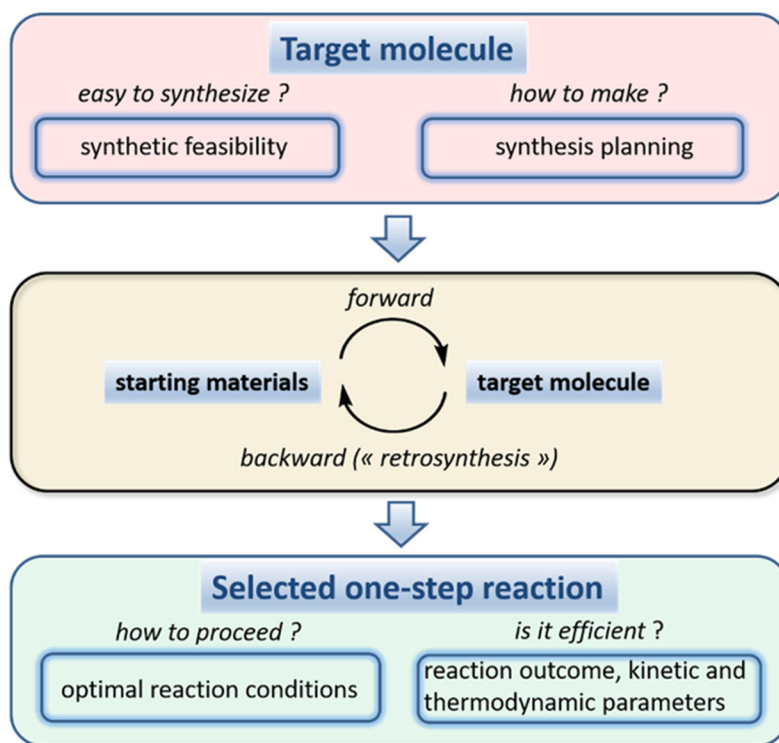


Figure 5. Main tasks of computer-aided synthesis design. As soon as a synthesis planning for a target molecule is established, efficiency of each one-step reaction and related optimal reaction conditions could be assessed.

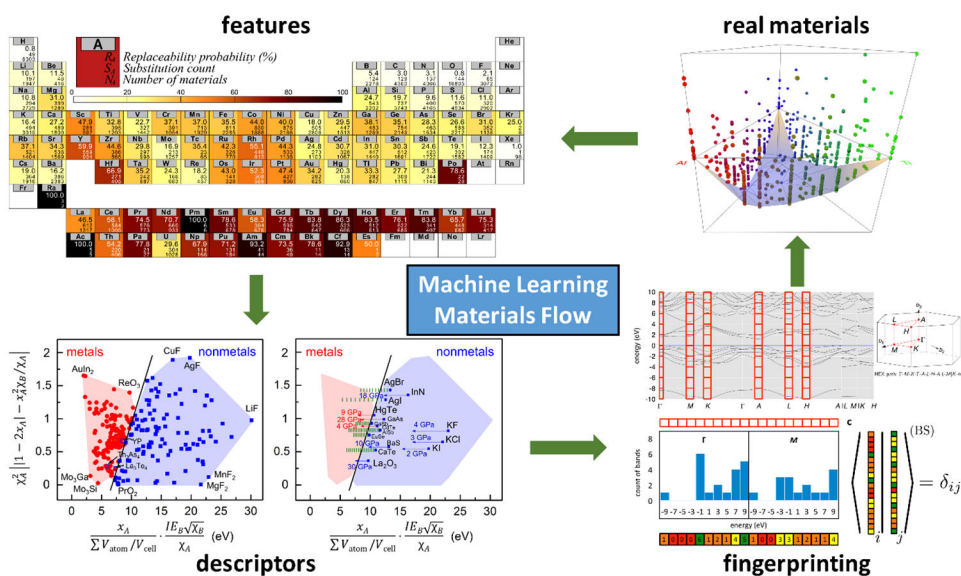


Figure 6. ML Materials Flow is a combination of feature extraction, descriptor analysis, structure fingerprinting (representations) of databases, and materials synthesizability. Figure reproduced from Refs.^{217,329-331}

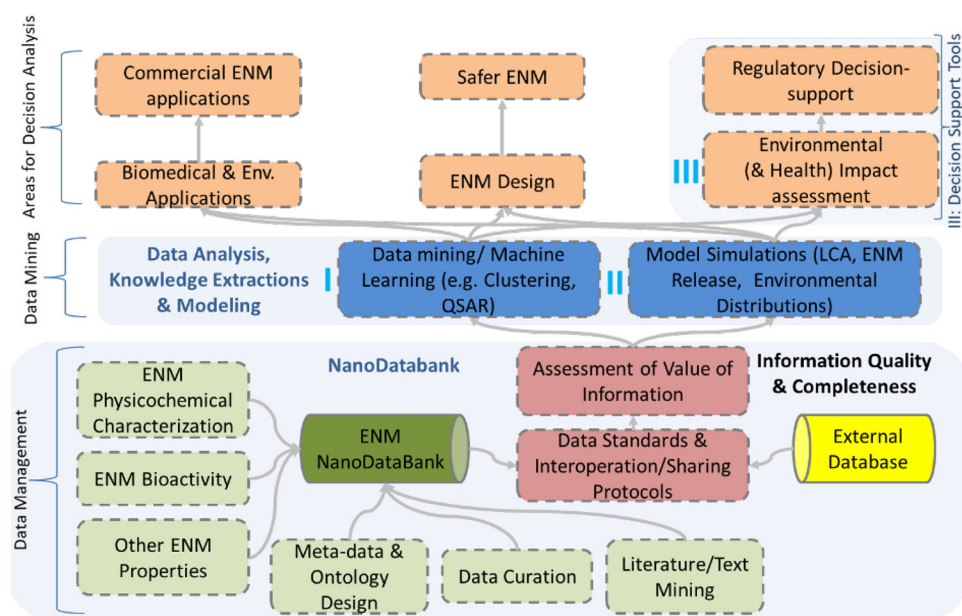


Figure 7. Nanoinformatics elements of environmental and health impact assessment for nanomaterials

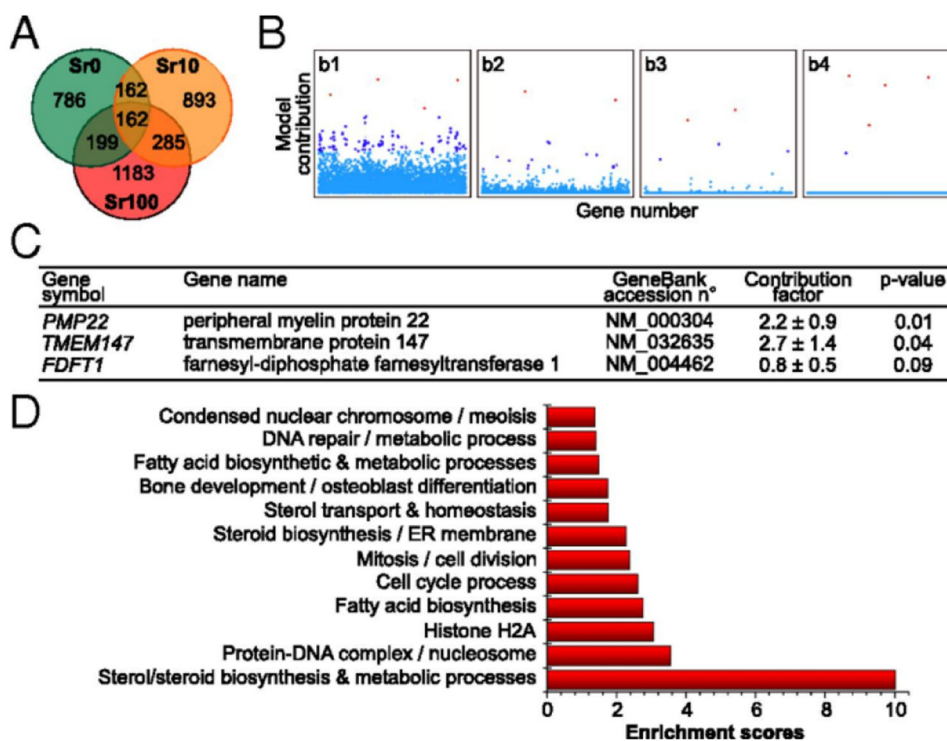


Figure 8. Changes in hMSC global mRNA expression mediated by treatment with BG- and SrBG-conditioned media. (A) Operation of the EM algorithm, showing progressive nulling of lower genes less relevant to the SrBG treatment. (B) The contribution (mean ± SE) of the most significant genes identified by sparse feature analysis. (C) Functional annotation clustering analysis of differentially expressed genes in response to Sr100 treatment compared with control. Reproduced from Ref.²⁵⁷



Figure 9.

SPE map of the correlation distances of the clinical and robotic parameters for the *completers* cohort. The map was derived by computing the pairwise Pearson correlation coefficients (R) for all pairs of features, converting them to correlation distances ($1-\text{abs}(R)$), and embedding the resulting matrix into 2 dimensions in such a way that the distances of the points on the map approximate as closely as possible the correlation distances of the respective features. The clinical parameters are highlighted in red, the robotic parameters on the affected side in blue, and the robotic parameters on the unaffected side in green. The map also shows distinct clusters of correlated variables which are preserved on both the affected and unaffected sides (outlined by green and blue ellipses, respectively).

Table 1.

Examples of QSAR-“inspired” studies from diverse research areas.

Cited Paper	Title	Journal	Year/Ref
292	Sensory analysis of red wines: Discrimination by adaptive fuzzy partition	Journal of Sensory Studies	2008/ ²⁹⁵
14	Improved wheat yield and production forecasting with a moisture stress index, AVHRR and MODIS data	Crop and Pasture Science	2009/ ²⁹⁶
14	Use of genetic algorithm and neural network approaches for risk factor selection: A case study of West Nile virus dynamics in an urban environment	Computers Environment and Urban Systems	2010/ ¹⁰
14	Whole cell-catalyzed transesterification of waste vegetable oil	Global Change Biology Bioenergy	2010/ ²⁹⁷
14	New Ground-Motion Prediction Equations Using Multi Expression Programing	Journal of Earthquake Engineering	2011/ ²⁹⁸
93	Qualitocracy: A Data Quality Collaborative Framework Applied to Citizen Science	IEEE Conference Proceedings	2012/ ²⁹⁹
14	Gene expression programming as a basis for new generation of electricity demand prediction models	Computers and Industrial Engineering	2014/ ³⁰⁰
292	Development of a model for quality evaluation of litchi fruit	Computers and Electronics in Agriculture	2014/ ³⁰¹
14,292	Good practices in LIBS analysis: Review and advices	Spectrochimica Acta Part B-Atomic Spectroscopy	2014/ ³⁰²
303	Characterization of Softwood and Hardwood LignoBoost Kraft Lignins with Emphasis on their Antioxidant Activity ¹⁵	BioResources	2014/ ³⁰⁴
292	Gene expression models for prediction of dam breach parameters	Journal of Hydroinformatics	2014/ ³⁰⁵
292	An entrainment model for non-uniform sediment	Earth Surface Processes and Landforms	2015/ ³⁰⁶
14	Indirect estimation of the ultimate bearing capacity of shallow foundations resting on rock masses	International Journal of Rock Mechanics and Mining Sciences	2015/ ³⁰⁷
14	A novel protocol for assessment of aboveground biomass in rangeland environments	Rangeland Journal	2015/ ³⁰⁸
14	Statistical Modeling of Soil Moisture, Integrating Satellite Remote-Sensing (SAR) and Ground-Based Data	Remote Sensing	2015/ ¹²
292	Testing and Prediction of Material Compatibility of Biofuel Candidates with Elastomeric Materials	International Journal of Fuels and Lubricants	2015/ ³⁰⁹
292	Regression Algorithms in Hyperspectral Data Analysis for Meat Quality Detection and Evaluation	Comprehensive Reviews in Food Science and Food Safety	2016/ ³¹⁰
292	Evolutionary patterns and physicochemical properties explain macroinvertebrate sensitivity to heavy metals	Ecological Applications	2016/ ³¹¹
292	Restricted attention to social cues in schizophrenia patients	European Archives of Psychiatry and Clinical Neuroscience	2016/ ³¹²
93	Molecular descriptor data explain market prices of a large commercial chemical compound library	Scientific Reports	2016/ ³¹³
14	A hybrid intelligent fuzzy predictive model with simulation for supplier evaluation and selection	Expert Systems with Applications	2016/ ³¹⁴
292	Development of a stage-dependent prognostic model to predict psychosis in ultra-high-risk patients seeking treatment for co-morbid psychiatric disorders	Psychological Medicine	2016/ ³¹⁵
292	Prediction of Timing of Watermain Failure Using Gene Expression Models	Water Resources Management	2016/ ³¹⁶
14	A new approach for modeling of flow number of asphalt mixtures	Archives of Civil and Mechanical Engineering	2017/ ³¹⁷
14	Next generation prediction model for daily solar radiation on horizontal surface using a hybrid neural network and simulated annealing method	Energy Conversion and Management	2017/ ³¹⁸

Cited Paper	Title	Journal	Year/Ref
93	Computer-Assisted Decision Support for Student Admissions Based on their Predicted Academic Performance	Journal of American Pharmaceutical Education	2017/ ¹¹
292	Predicting Bond Strength between FRP Plates and Concrete Substrate: Applications of GMDH and MNLN Approaches	Journal of Advanced Concrete Technology	2017/ ³¹⁹
14	Gene Expression Programming Approach to Cost Estimation Formulation for Utility Projects	Journal of Civil Engineering and Management	2017/ ³²⁰
292	Prediction of flow duration curves for ungauged basins	Journal of Hydrology	2017/ ³²¹
14	Maize [<i>Zea Mays</i> (L.)] crop-nutrient response functions extrapolation for Sub-Saharan Africa	Nutrient Cycling in Agroecosystems	2017/ ³²²
14	Performance assessment of existing models to predict brittle failure modes of steel-to-timber connections loaded parallel-to-grain with dowel-type fasteners	Engineering Structures	2018/ ³²³
292	A comparative study on groundwater spring potential analysis based on statistical index, index of entropy and certainty factors models	Geocarto International	2018/ ³²⁴
325	Environmental factors influencing snowfall and snowfall prediction in the Tianshan Mountains, Northwest China	Journal of Arid Land	2018/ ³²⁶
14,292	Prediction of riprap stone size under overtopping flow using data-driven models	International Journal of River Basin Management	2018/ ³²⁷
14	Forecasting experiments of a dynamical-statistical model of the sea surface temperature anomaly field based on the improved self-memorization principle	Ocean Science	2018/ ⁹
292	Expressed emotion as a predictor of the first psychotic episode - Results of the European prediction of psychosis study	Schizophrenia Research	2018/ ³²⁸