# Artificial intelligence in early drug discovery enabling precision medicine

Fabio Boniolo, Emilio Dorigatti, Alexander J. Ohnmacht, Dieter Saur, Benjamin Schubert & Michael P. Menden

Published online: 02 Jun 2021.

Submit your article to this journal ⬀

Article views: 4390

View related articles ⬀

View Crossmark data ⬀

Citing articles: 1 View citing articles ⬀

Taylor & Francis
Taylor & Francis Group

REVIEW

# Artificial intelligence in early drug discovery enabling precision medicine

Fabio Boniolo[a,b,*], Emilio Dorigatti[a,c,*], Alexander J. Ohnmacht[a,d,*], Dieter Saur[b], Benjamin Schubert[a,e] and Michael P. Menden [a,d,f]

aInstitute of Computational Biology, Helmholtz Zentrum München - German Research Centre for Environmental Health, Munich, Germany; bSchool of Medicine, Chair of Translational Cancer Research and Institute for Experimental Cancer Therapy, Klinikum Rechts Der Isar, Technische Universität München, Munich, Germany; cStatistical Learning and Data Science, Department of Statistics, Ludwig Maximilian Universität München, Munich, Germany; dDepartment of Biology, Ludwig-Maximilians University Munich, Martinsried, Germany; eDepartment of Mathematics, Technical University of Munich, Garching, Germany; fGerman Centre for Diabetes Research (DZD e.V.), Neuherberg, Germany

**ABSTRACT**

**Introduction**: Precision medicine is the concept of treating diseases based on environmental factors, lifestyles, and molecular profiles of patients. This approach has been found to increase success rates of clinical trials and accelerate drug approvals. However, current precision medicine applications in early drug discovery use only a handful of molecular biomarkers to make decisions, whilst clinics gear up to capture the full molecular landscape of patients in the near future. This deep multi-omics characterization demands new analysis strategies to identify appropriate treatment regimens, which we envision will be pioneered by artificial intelligence.

**Areas covered**: In this review, the authors discuss the current state of drug discovery in precision medicine and present our vision of how artificial intelligence will impact biomarker discovery and drug design.

**Expert opinion**: Precision medicine is expected to revolutionize modern medicine; however, its traditional form is focusing on a few biomarkers, thus not equipped to leverage the full power of molecular landscapes. For learning how the development of drugs can be tailored to the heterogeneity of patients across their molecular profiles, artificial intelligence algorithms are the next frontier in precision medicine and will enable a fully personalized approach in drug design, and thus ultimately impacting clinical practice.

## 1. Introduction

It is estimated that the cost of discovering and developing a new drug is around USD 3 billion [1], with an approval rate close to 13% when considering all the compounds that reach clinical trials [2]. In particular, novel compounds register an underwhelming success rate of ~66% in Phase I regarding tolerability and side effects, ~48% in Phase II concerning dosage and efficacy, and ~59% in Phase III regarding efficacy and toxicology [2].

Precision medicine has raised high hopes to improve the success rate of Phase II and III clinical trials by tailoring treatment options to the characteristics of patient subgroups based on differences in molecular profiles, lifestyle, and environmental factors. Since the introduction of this paradigm, the number of application areas in medicine and healthcare has rapidly increased, with oncology being the vanguard for its deployment [3]. Beyond that, recent works on cardiovascular diseases [4,5], type 2 diabetes [6] and neurodegenerative disorders, such as Alzheimer's disease [7] and Amyotrophic Lateral Sclerosis [8,9] have highlighted the growing relevance of precision medicine on the whole healthcare sector.

Large leaps forward in precision medicine were achieved by the rapid development of DNA-sequencing technologies and their regular use in clinical practice [10]. Since then the deep molecular characterization of patients was expanded to transcriptomics [11], epigenomics [12], and proteomics [13], collectively referred to as 'omics' technologies. These technical developments, together with advances in information technology, computer science, and computational biology, have created a fertile ground for the successful integration of artificial intelligence (AI) with precision medicine.

Conventional drug development pipelines consist of target identification and validation, assay development and screening, hit identification, lead optimization, and the selection of the final molecule for clinical development [14], each step marking a milestone in a rigid streamlined process. Main objectives are to identify potent drugs with suitable bioavailability, toxicity profiles, chemical synthesis, selectivity against putative target and ADME (absorption, distribution, metabolism, excretion), whilst mostly neglecting the heterogeneity of patients. In order to address this, precision medicine was introduced to customize treatments based on patient profiles [15]. This concept strongly impacted the linearity of drug

---

CONTACT Benjamin Schubert ✉ benjamin.schubert@helmholtz-muenchen.de; Michael P. Menden michael.menden@helmholtz-muenchen.de Institute of Computational Biology, Helmholtz Zentrum München - German Research Centre for Environmental Health, 85764 Munich, Germany

*equal contribution

discovery pipelines and suggested a more integrated and looped process [16]. Strong benefits of response biomarker discovery include acceleration of drug approval due to increased success rates of clinical trials [2] and thus reduced costs. AI offers the potential to leverage the entire molecular landscape of patients, thus becoming an invaluable tool for precision medicine.

The identification of actionable disease subtypes is increasingly powered by AI using integrative methods combining diverse data modalities [17–26]. Complementarily, pharmacogenomics benefits from machine learning methods predicting *in vitro* monotherapy response by using molecular data in cell cultures, which may yield novel predictive biomarkers through the interpretation of the learned models [27–34]. Likewise, AI showed encouraging progress in the discovery of potent drug combinations [35–39].

A concept which is often utilized in conjunction with precision medicine is drug repurposing, which leverages drugs approved for one disease indication and applies them in another context. For example, the discovery of EML4-ALK fusions in non-small cell lung carcinoma (NSCLC) patients triggered the repurposing of crizotinib [40], which is a potent ALK, MET and ROS1 inhibitor and was not developed for this indication in the first place. Remarkably, crizotinib was approved within a record time of only four years within this new indication, making this an exemplar of both drug repurposing and precision medicine. Nowadays, AI is increasingly used for drug repurposing by systematically aggregating various omics layers and drug features for training models that prioritize compounds based on their properties [41,42].

Personalized medicine is the extreme case of precision medicine, in which the treatment is not only administered according to a biomarker but truly tailored to the needs of an individual. AI has been successfully used to develop individualized drug compounds themselves, with personalized cancer vaccines being one of the prime examples [43–47]. Cancer vaccines require the identification of antigen peptides that are highly specific to the patient's tumor and MHC genotype and use those to boost the patient's immune system [48]. Machine learning [49–51] and optimization methods have been developed to aid peptide identification and assembly of the vaccine

[52–54] and have been integrated in almost all personalized vaccine design pipelines. The ability to choose a target antigen and set of MHC alleles makes such vaccine design frameworks not only applicable to personalized cancer immunotherapy [43–47], but generally useful for population-level prophylactic vaccine development against infectious diseases. Similarly, large and small molecule design has also seen recent successes [55–57] in AI-driven drug development and even some examples of personalized applications [58].

In spite of these early successes, however, such AI-powered approaches still need to be translated into standard clinical practice, and a fully *in silico* drug design approach that integrates personalized patient information has yet to be realized. Nonetheless, AI has already impacted multiple sectors of the pharmaceutical industry [59] and the recent advances suggest that its applications may enable precision medicine in the clinics in the near future.

Thus, this review focuses on the impact of AI on drug discovery and development by building a bridge between biomarker discovery and drug design illustrated through its pioneering applications in precision oncology. We will particularly focus on the role of such techniques in biomarker discovery via disease subtyping (Figure 1(a)), high-throughput screens (Figure 1(b)) and drug combinations (Figure 1(c)), as well as in drug design (Figure 1(d-f)).

## 1.1. Artificial intelligence

Applications of AI led to radical changes both in academia and in industry, often disrupting the typical approaches to science and business while introducing new methodological, epistemological [60], ethical [61], and privacy-related [62] concepts. In particular, AI is becoming increasingly important in biomedicine and healthcare [63], where it led to breakthroughs in biomedical image analysis [64], prognosis [65], patient care [66], and clinical decision support [67]. This review focuses on AI in early drug discovery enabling precision medicine (Table 1).

Over the years, a wide variety of algorithms has been applied to a multitude of predictive tasks in precision medicine. Generally, predictive power and interpretability of the algorithm's decisions are inversely related.

One contemporary research area in AI is focusing on enabling machines to learn general rules from provided example data, also known as machine learning, to make predictions for previously unseen samples [68]. Machine learning can be further categorized into supervised learning, in which phenotypic observations are known and relations between input features and these observations are sought. In contrast, unsupervised learning aims to uncover hidden patterns in the data by clustering or latent factor modeling to explain observed variability. Alongside these two fundamental paradigms of machine learning, reinforcement learning gains more traction in biotechnological sciences, especially in drug development. Here, the model is allowed to take actions, such as introducing an amino acid alteration, in a pre-specified environment (a protein) to optimize a specific property (efficacy).

A wide spectrum of learning algorithms has been proposed. They differ in the complexity of concepts they are

**Figure 1.** AI enhances biomarker discovery and drug design. (a) Data-driven disease subtyping gives insights into the diverse disease etiologies, which can be leveraged for patient stratification. (b) Functional genomics and drug screens systematically explore large panels of disease models, facilitating the discovery of novel biomarkers with machine learning. (c) AI algorithms guide the prioritization of potent drug combinations to overcome monotherapy resistance. (d) AI accelerates and guides vaccine aiding in each step of the design process. Through deep generative models the design space of (e) proteins and (f) small molecules can be systematically searched to generate novel drugs that are difficult to attain through experimental designs.

**Table 1.** Overview of AI algorithms used in precision medicine.

| AI Algorithm | | Advantages | Disadvantages | Applications discussed (section) |
|---|---|---|---|---|
| Shallow Learning | Linear/Logistic Regression | + Interpretable | - Limited to linear trends | * drug response (2.2)<br>* drug combinations (2.3)<br>* MHC affinity (3.1)<br>* T-cell specificity (3.1) |
| | Support Vector Machines (SVM) | + Nonlinear function approximation<br>+ Flexible through kernel | - Less interpretable<br>- Hard to design kernels for nonstandard data | * MHC affinity (3.1) |
| | Random Forests | + Nonlinear function approximation<br>+ Automatic handling of different data types<br>+ Interpretable | - Not well equipped for regression tasks<br>- Less interpretable | * disease subtyping (2.1)<br>* patient stratification (2.1)<br>* drug response (2.2)<br>* drug combinations (2.3)<br>* T-cell specificity (3.1) |
| | Gaussian Processes | + Nonlinear function approximation<br>+ Flexible through kernel<br>+ Fully Bayesian | - Does not scale well to large datasets<br>- Hard to design kernels for nonstandard data | * MHC affinity (3.1)<br>* T-cell specificity (3.1) |
| | Dimensionality reduction and feature synthesis | + No labels needed | - Limited expressive power | * disease subtyping (2.1)<br>* patient stratification (2.1)<br>* drug response (2.2) |
| Deep Learning | Generative | + Nonlinear<br>+ Scales well<br>+ Handles unstructured data<br>+ Can generate novel examples<br>+ Few labels are needed (if at all) | - Hard to interpret<br>- Needs lots of data and compute resources<br>- Novel examples can be hard to evaluate | * protein sequence (3.2) modeling<br>* small molecule modeling (3.3) |
| | Discriminative | + Nonlinear<br>+ Scales well<br>+ Handles unstructured data | - Hard to interpret<br>- Needs lots of data and compute resources | * disease subtyping (2.1)<br>* drug response (2.2)<br>* drug combinations (2.3)<br>* MHC affinity (3.1)<br>* T-cell specificity (3.1) |

able to learn, the type of examples they are able to handle, and to which degree they can provide comprehensible explanations for their decisions. Linear models are easily interpretable but limited to very simple relationships between features and concepts. Random forests are an ensemble of models that apply a sequence of learned rules to reach a conclusion via majority voting [69]. While these rules remain easily interpretable, the decision boundary of random forests can be quite complex and nonlinear. Support vector machines use flexible user-defined feature extractors to identify training examples that are similar to the provided query in a new transformed latent space [70].

In contrast to most other machine learning methods, deep neural networks are able to automatically extract highly complex patterns from all sorts of data types [71,72]. Data-hungry and often inscrutable in how a deep learning network comes to its conclusion, its predictions can be in many cases vastly more accurate than other methods, in particular in unstructured domains, such as images or molecular entities, when applied to very large datasets and tuned appropriately [73], but at the same time, they are highly nonlinear and thus more challenging to interpret [74]. To overcome this, a subfield has emerged, called explainable AI, that studies and develops methods to gain a better understanding of how AI algorithms come to their conclusions. Such approaches found large appeal in many high-risk application areas such as precision medicine [75]. This intriguing field, however, is still in its infancy especially for complex models, such as deep neural networks.

## 2. Artificial intelligence in biomarker discovery

The identification of robust disease-specific biomarkers is fundamental for advancing precision medicine. Therefore, it is of paramount importance to develop methods able to identify actionable molecular targets suitable for therapy [76]. State-of-the-art precision oncology mostly focuses on well-characterized cancer somatic mutations, and in some instances on germline variants, for patient stratification [77]. AI offers opportunities for identifying complex biomarker signatures, i.e. disease-specific altered networks of genes and proteins shared between tumor types and across multi-omics layers that overcome the obsolete 'one gene, one drug, one disease' paradigm [78].

### 2.1. Disease subtyping

Disease subtyping is a powerful concept to reduce the dimensionality of complex disease characteristics into simplified signatures, which can be used for patient stratification. Particularly, AI methods have pioneered the development of robust and clinically relevant disease subtypes [79]. The concept of disease subtyping has been leveraged in treating many diseases [5,7], however, oncology is the prime example for this.

It is important to understand that cancer is not a unique disease but rather a heterogeneous category of diseases. Tumors are usually characterized by their tissue of origin, and further classified into molecular tumor subtypes [80]. In addition, tumor subtypes may be defined based on clinical information [81] and imaging [82]. For example, tumors showing a higher rate of lesion enhancement on MR images are more likely classified as luminal B subtype in breast cancer and are independently associated with better prognosis of patients [83].
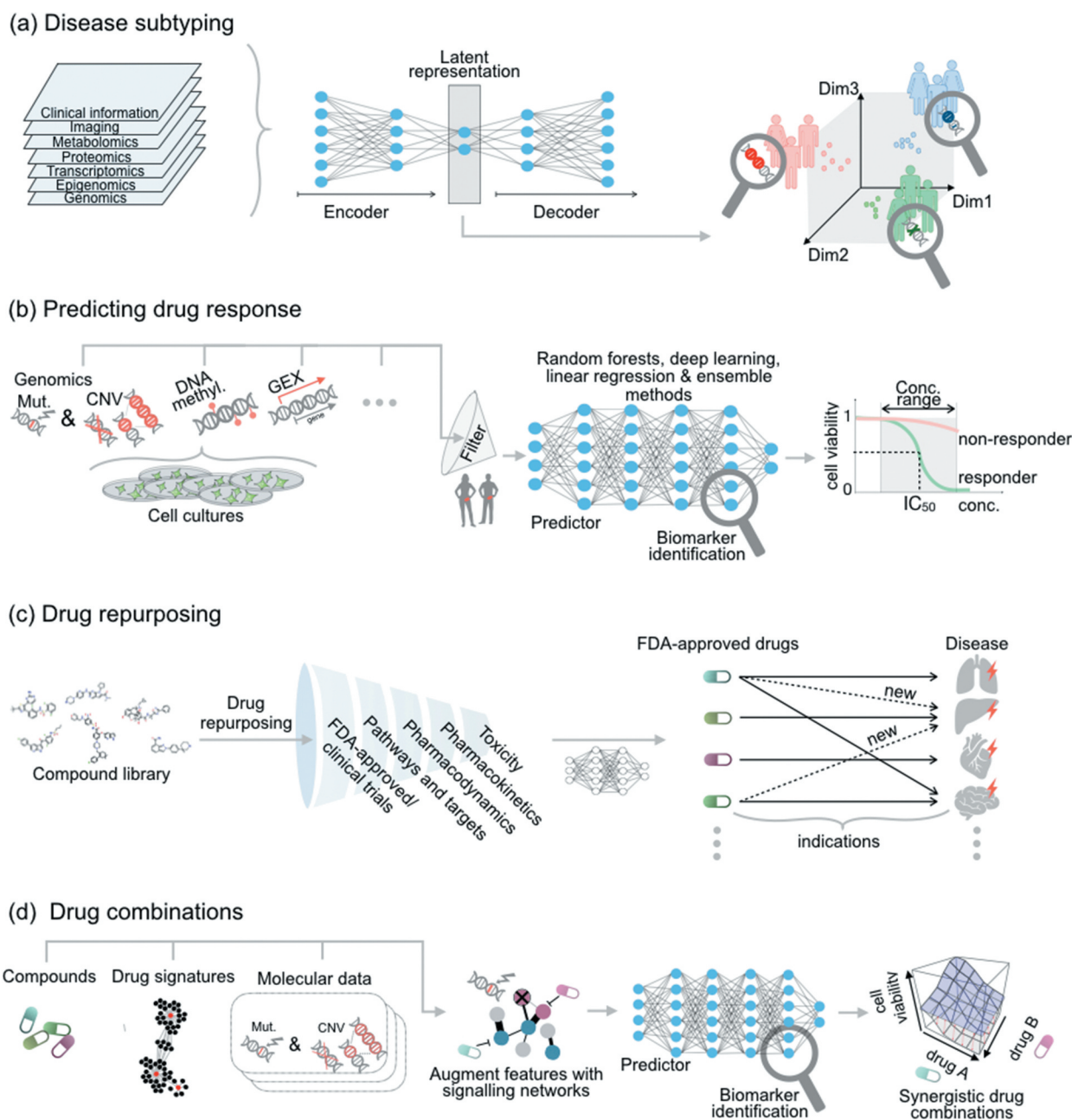
Colorectal cancer research, which often lacks univariate genetic biomarkers of drug sensitivity [84], pioneered gene expression subtyping efforts utilizing machine learning, and concluded that there are four consensus molecular subtypes [17]. This effort, based on random forest classification, was led by six leading colorectal cancer research teams who defined the optimal number of biologically relevant clusters that ultimately proved to have prognostic value in clinical trials [85].

Following the example of colorectal cancer subtyping, other cancer entities without clear genetic biomarker indications successfully used gene expression subtyping, e.g. in bladder cancer [18] and pancreatic cancer [86]. The main characteristics of cancer subtypes are tumor differentiation stages: well-differentiated (classical) subtypes seem to be less aggressive and more sensitive to therapies compared to undifferentiated (basal) tumors [87]. These efforts used unsupervised machine learning techniques, i.e. non-negative matrix factorization (NMF) and independent component analysis (ICA), and consecutively, determined the optimal number of signatures according to biological interpretability and clinical translatability.

The rise of large-scale molecular datasets, such as the Cancer Genome Atlas (TCGA) [88], the International Cancer Genome Consortium (ICGC) [89], and the Pan-Cancer Analysis of Whole Genomes (PCAWG) [90], empowered deep learning approaches to systematically derive robust disease subtypes. For example, variational autoencoders have proven to be successful in identifying targetable subtypes of non-small cell lung cancer when trained on methylation data [19]. Similar approaches using different omics layers and various architectures, such as mixtures of unsupervised and supervised components, led to the stratification of patients in neuroblastoma [20], lung adenocarcinoma [21], and breast cancer [22].

Most tumor subtyping efforts solely explore one data modality. Therefore, they often neglect complementary information and dependencies contained in different molecular layers [91]. To address this, attempts have been directed to the development of machine learning techniques in multi-omics integration [23]. These methods jointly model data from various molecular sources [92] and extract directions of common variance using, e.g. canonical correlation analysis (CCA) [24] or factor analysis [25]. Emerging AI-based methods such as variational autoencoders offer a chance to further improve data integration by allowing the projection of multiple omics layers [26], or even of different data modalities [93], to the same latent space (Figure 2(a)).

In summary, disease subtyping is driven by unsupervised learning techniques exploring deep molecular landscapes, and in conjunction with expert domain knowledge, may deliver interpretable and actionable biomarkers. This is complementary to the state-of-the-art precision medicine approach, which mostly focuses on 'one-gene' associations. We envision that the identification of disease subtypes with AI and their

**Figure 2.** AI in early drug and biomarker discovery. (a) Deep learning empowers precision medicine and disease subtyping by revealing meaningful patient subgroups based on molecular and clinical data. (b) High-throughput drug screens in cell cultures, in conjunction with deep molecular characterization of these cell cultures, are leveraged to predict drug response and identify biomarkers. (c) Drug repurposing identifies new therapeutic applications of existing drugs. (d) Predicting drug synergies guides the prioritization of synergistic drug combinations for increased treatment efficacies.

interaction with established single-gene biomarkers may drive patient stratification, and thus the design of subtype-specific treatment options in the near future.

## 2.2. High-throughput screens

Preclinical *in vitro* and *in vivo* disease models are ranging from simple immortalized human cancer cell lines [40,94,95], self-organizing 3D cell cultures, i.e. organoids [96], to complex animal models [97]. Biological models with increased complexity more accurately capture human tumor biology, whilst simpler models enhance a systematic and parallel comparison of many therapeutic agents across many samples.

In particular, high-throughput drug screens (HTS) applied to *in vitro* 2D or 3D cell cultures provide a rich resource to study pharmacogenomic interactions when complemented with deep molecular characteristics. Cell cultures remain simplified models which are accompanied by culturing artifacts, mostly disregarding the tumor microenvironment and immune responses. Nevertheless, these simplistic models recapitulate clinically relevant aspects [98,99]. Such pharmacogenomics datasets are well suited for the application of machine learning methods for predicting drug sensitivity and identifying drug response biomarkers [100].

Drug response in cancer HTSs is typically quantified either by the drug concentration required to reduce cell viability by half

(IC$_{50}$ value), or by the area under the dose-response curve (AUC). In contrast, functional genomics approaches quantify gene essentiality based on pooled genome-wide CRISPR-Cas9 activation or depletion screens, e.g. in the Cancer Dependency Map (DepMap) project [101]. Such HTSs enable the analysis of genetic dependencies in cancer, leading to the identification of cancer vulnerabilities and thereby revealing potential drug targets [102].

Experimental HTS technologies spawned a multitude of novel computational methods to predict drug response leveraging the multi-omics characterization of cell lines [27,103] (Figure 2(b)). Transfer learning approaches have been used to enable prediction in a single cancer type with a small number of samples by leveraging information about related cancer types with many samples [104]. In addition, multi-task learning methods make multiple drug sensitivity predictions simultaneously by integrating cell line and drug features, i.e. quantitative structure–activity relationship (QSAR) models [28]. Likewise, similarity patterns in drug response across cell lines and drugs have been quantified with constrained matrix factorization [29], or by projecting high-dimensional feature spaces into low-dimensional embedding spaces that correlate better with drug-relevant pathways [30]. In addition, deep learning methods are leveraged for analyzing high-throughput screens, e.g. convolutional architectures show promising results [31] and autoencoders are progressively used for these tasks [32]. These models are able to prioritize novel compounds for HTS, and some approaches can identify biomarkers when exploring their feature weights [33]. For example, a visible neural network trained on tumor genotypes and chemical structures has been shown to reach good predictive performances while identifying pathways involved in the response to mTOR inhibitors (e.g. everolimus) or CDK4/6 inhibitors (e.g. ribociclib) in ER positive metastatic breast cancer [34].

In contrast to conventional drug HTSs, which explore baseline molecular characteristics to predict drug response, the connectivity map (C-Map) focused on deriving drug signatures by measuring gene expression or proteomics before and after treatment [105,106]. This HTS concept enables drug repurposing [107], a strategy to prioritize compounds that are already approved in one disease indication, and being reused within another indication. As a consequence of drug repurposing, it is possible to exploit the already existing toxicity, pharmacokinetics, and pharmacodynamic profiles, thus accelerating clinical development. Typical AI approaches for drug repurposing involve the use of databases such as ChEMBL [108] and PubChem [109] containing bioactivity profiles of compounds and ADMET (absorption, distribution, metabolism, excretion, toxicity) containing compound properties. By leveraging these databases, together with interaction networks and molecular data, machine learning methods have been developed for drug target prediction [110–113] and drug repurposing [114–117], including deep learning methods [118,119].

The availability of extensive molecular data repositories and the advances in computational analysis approaches offer the opportunity to drive a more systematic methodology, thus combining drug repurposing, AI and precision medicine (Figure 2(c)) [120]. The current COVID-19 pandemic has demanded new treatment options for the disease and has spawned novel AI strategies for prioritizing candidate drugs via drug repurposing for their accelerated usage [121], i.e. exploiting relationships of protein–protein interactions between drug targets and SARS-CoV-2 viral targets, thus short-listing 81 candidate drugs for COVID-19 treatment [121].

In conclusion, the computational analysis of high-throughput functional genomics and drug screens offers a route to identify drug targets and response biomarkers by leveraging different layers of molecular and chemical information. These methodologies yet have to find their way into the routine practise, but the increasing size of available datasets, both the deep molecular characterization of biological models and treatment patterns, in conjunction with advances in machine learning, enables the development of effective strategies for drug repurposing and precision medicine.

## 2.3. Drug combinations

Monotherapies can suffer from low potency, in particular, acquired drug resistance is a major obstacle in oncology [122,123]. In order to address this, drug combinations may exploit drug synergies (Figure 2(d)) to anticipate tumor evolution and overcome resistance [37]. Most drug combinations follow either a strategy to 'double-hit' the same signaling pathway, or alternatively, they target independent pathways or disease mechanisms [37,122].

The space of possible drug combinations grows exponentially when exploring drug cocktails, since a set of n drugs can form $2^n - 1$ unique subsets, thus highlighting the need for computational methods to prioritize the most promising combinations, whilst not increasing toxicity. In recent years, drug combinations HTSs were established [37,124], thereby creating a fertile ground for novel machine learning models that predict drug synergy.

Traditionally, methods for discovering potent drug combinations predominantly used systems biology approaches and were only using drug combination data as validation experiments. For example, drug signatures were leveraged to derive drug functional networks, in which potent drug combinations were extracted by searching drugs whose targets were enriched in a complementary disease specific signaling network [35]. In a broader context, methods imposed similarity metrics between drug and disease signatures, following the paradigm that an ideal drug combination would fully reverse a given disease signature [36].

In the context of the AstraZeneca-DREAM crowd-sourcing challenge [37], drug synergy prediction methods were able to leverage the molecular landscape of cancer cell lines (Figure 2 (d)). The best performer used random forests and incorporated prior knowledge of protein–protein interactions (PPI) and cancer signaling networks to augment their feature set. Competing methods based on deep learning frameworks were among the top three performers, and recently, other AI-based models gained more traction in this field [38,39].

Despite the large number of computational approaches for predicting drug combinations, we are still lacking synergy biomarkers in clinics. The bottlenecks are the sizes of currently

available drug combination HTSs, distinguishing between *in vitro* drug synergy and toxicity, and experimental *in vivo* validations. Furthermore, such methods usually focus solely on synergy and neglect acquired resistances and adverse effects [125]. However, future efforts will expand the dimensions of drug combination HTS in conjunction with *in vivo* validation, thereby driving the next generation of prioritization algorithms of drug combinations.

## 3. Artificial intelligence in drug design

AI methods have a long-standing tradition in drug development, and the combination of high-throughput experimental techniques with deep learning has led to some impressive results that we are now starting to see in both precision and non-precision drug development [100,126]. In the conventional linear drug discovery pipeline, drug design follows target identification and is validated through clinical trials, most of which fail due to lack of efficacy or unwanted side-effects [2]. Not only can biomarkers improve the success rate of clinical trials by identifying likely responders for inclusion in the trial [127], they can also provide novel mechanistic insights on basic biology and disease pathogenesis, guiding the development of novel drug treatments with improved immunogenicity and reduced toxicity, targeted for example to disease-causing genetic variants, mutant proteins, or their associated pathways [15]. The AI-based drug design methods discussed in this section could in the future be adapted to disease-specific characteristics by tuning the objective function that is being optimized and evaluating them against more appropriate benchmarks.

### 3.1. Vaccine design

Cancer immunotherapies exploit the patient's immune system to fight tumors [128]. During tumor evolution, unique genomic alterations might arise that give rise to neoepitopes – a class of major histocompatibility complex (MHC)-bound, altered self-peptides that differentiate healthy and tumoral cells. Such cancer-specific peptides have been successfully exploited by AI-based methods for personalized vaccine design [43–47], and since then AI has become an indispensable component of many cancer vaccine design pipelines (Figure 3(a)) – from predicting neoepitopes from a patient's uniquely altered peptide pool to selecting and assembling the neoepitopes into vaccines.

To characterize the tumor surface and identify its neoepitopes, several predictive models have been proposed including kernel methods [49,129], position-specific scoring matrices [50,130], and neural networks [51,131,132]. Advances in deep learning and mass spectrometry [133] have made it possible to train increasingly complex models on larger datasets [134,135] and accounting for additional information (e.g. tumor antigen expression), which has led to superior performance [136]. These prediction tools have now been bundled in software pipelines that can generate predicted neoepitopes from a provided list of somatic mutations [137–139].

To predict immune reactive neoepitopes, a recent consortium effort has benchmarked a large variety of different features, including the MHC-binding affinity ratio between neoepitope and wildtype peptide, similarity to viral peptide sequences, and physicochemical properties, presumed to be associated with neoepitope immunogenicity. The study found that most T-cell reactive neoepitopes were highly foreign or agretopic [140]. Others have started to directly model the T-cell-epitope interaction using linear models [141,142], Gaussian processes [143], random forests [144], and deep neural networks [145,146] but all with limited generalization capability to unseen epitope-T-cell pairs due to data limitations.

In cancer vaccine design, tens to thousands of possible neoepitopes are initially identified depending on the cancer type [147] using predictive models. Subsequently, selecting a small set of neoepitopes that is maximally immunogenic with respect to the MHC alleles of the patient or target population, while at the same time guaranteeing sufficient diversity to cover tumor or pathogen heterogeneity is of prime importance in producing an effective vaccine. This discrete optimization problem was first approached by *ad-hoc* scoring rules that consider the neoepitopes' antigen coverage and other relevant qualities [54,148]. Successive approaches used integer linear programming to optimally solve the selection problem [53]. These methods were complemented by other computational techniques addressing the assembly of the selected neoepitopes into vaccines and optimizing their intracellular processing [52,149], which were recently experimentally shown to improve vaccine efficacy over manual designs [150]. Eventually, these successive approaches were unified into a single vaccine design framework that allows modeling the trade-off between immunogenicity and vaccine processing [151,152].
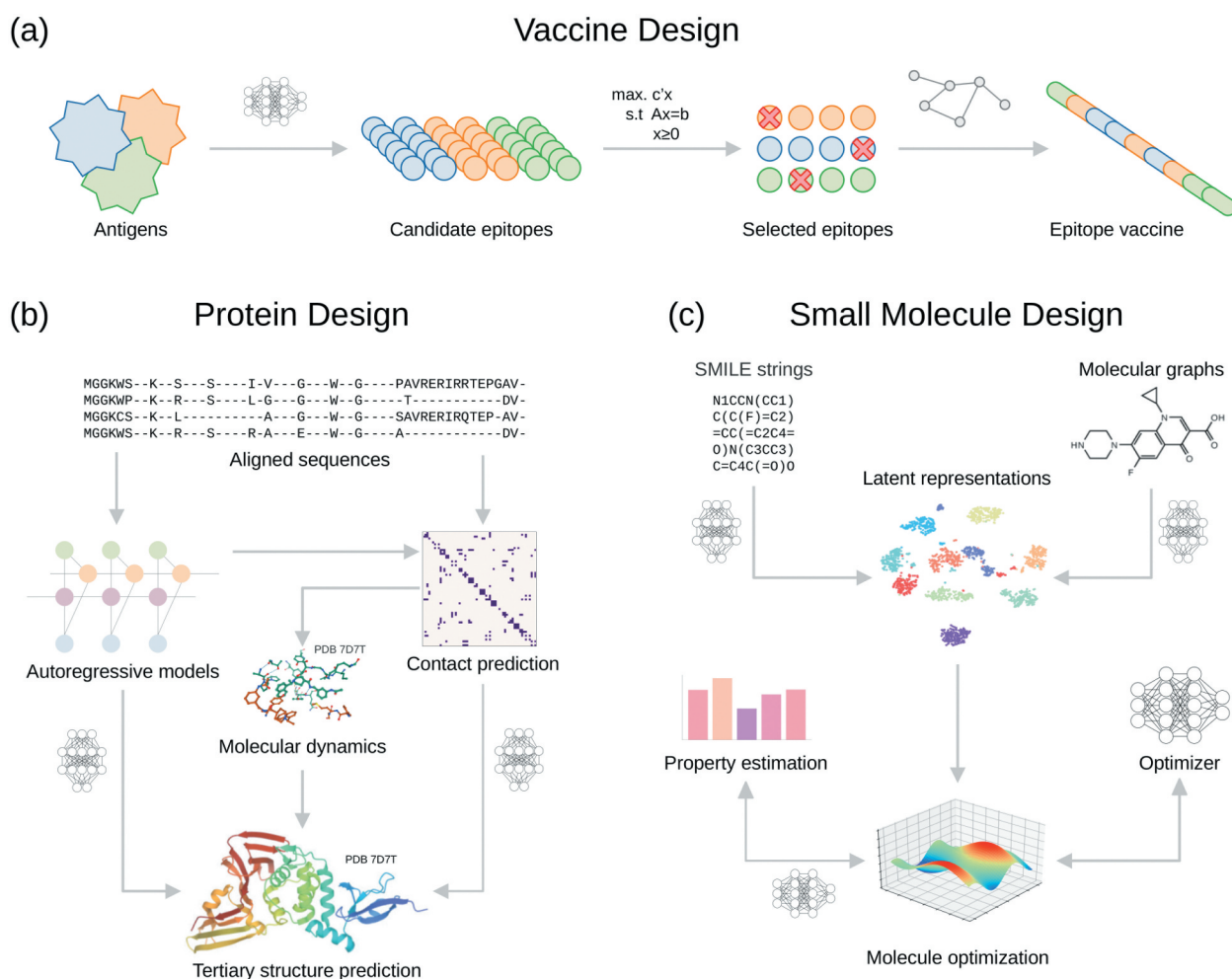
The first clinical applications of AI-driven personalized immunotherapies were quite encouraging, yet the objective functions used to select neoepitopes in vaccine design are still arbitrary and based on presumptions. A deeper understanding of what neoepitopes to target and which metrics lead to good vaccine formulas is necessary. To this end, large concerted experimental screening efforts are required. Similarly, their routine clinical application necessitates easy-to-use interfaces and computational analysis pipelines seamlessly integrated into the standard sample processing and diagnostic analyses.

In summary, AI advanced the identification and ranking of neoepitopes arising from patients' tumors, as well as the composition and formulation of epitope vaccines for efficient processing and maximal effectiveness. Such approaches are however general enough and can be used with little modifications to design prophylactic vaccines optimized for a target population and pathogen pool based on MHC and strain prevalence [53], examples of which are HIV [153], Influenza [154], Malaria [155], Hepatitis C virus [156], and SARS-CoV -2 [157].

### 3.2. Protein design

Protein engineering has seen a recent surge of innovations driven by breakthroughs in unsupervised representation learning of protein sequences, efficiently exploiting hundreds of millions of sequenced proteins to learn which positions co-evolve in a large set of evolutionary diverse dataset of the

**Figure 3.** AI approaches in drug development. (a) AI is used to aid in each step of modern cancer vaccine pipelines. Supervised prediction models are used to identify cancer-specific antigenic peptides and constraint optimization models are employed to select and assemble the final vaccine. (b) Generative models that use co-evolutionary information of protein sequences have revolutionized structural biology in recent years and are now applied for 3D structure prediction, accelerate molecular dynamics simulations, and to design novel proteins. (c) Deep neural networks are currently revolutionizing the design process of novel small molecules. With generative models, novel molecular structures with optimized biochemical properties can easily be created exploring large proportions of the chemical search space.

same protein family (Figure 3(b)). Through such co-evolutionary information, AI can learn which residues are critical for the function of the protein and which residues can be altered to tailor its properties, making it possible to adapt existing therapeutics for specific patients while reducing potential immune-related side effects [58].

Co-evolutionary protein sequence modeling has been shown to be predictive for the 3D structure of proteins [158,159] and complexes [160–162], but also for mutation effects [163,164]. Early approaches were based on maximum entropy models inspired by statistical physics but have since been superseded by deep learning-based models, such as variational autoencoders [164,165], generative adversarial networks [166], autoregressive models borrowed from natural language processing [167–171] and novel architectures adapted for the task [172].

The use of co-evolutionary sequence information has also led to new breakthroughs in protein structure prediction, as demonstrated in the two latest CASP competitions by AlphaFold [173]. This model uses predicted structural contacts

from co-evolutionary models and refines them with a deep residual network that predicts the distribution of contact distances. These are interpreted as the statistical energy of the protein fold and directly minimized to yield highly accurate 3D protein structures. Since then, extensions have been made to remove the dependency on co-evolutionary analysis further [174], making it possible to predict 3D structures of artificial proteins and protein families starting from only a few sequences. Others have started to develop neural protein folding simulators that are fully trainable in an end-to-end-fashion and can directly map sequence to structure [175,176].

Such generative models are now being slowly adopted for protein engineering. One of the earliest approaches used maximum entropy models in combination with integer linear programming to re-engineer existing biotherapeutics with reduced adverse immune reactions, focusing on editing immunogenic regions of the biotherapeutic based on pre-specified MHC molecules and opening up the possibility for personalized re-engineered biotherapeutics [58]. Others used deep autoregressive models to optimize screening libraries of

nanobodies for higher functional yield [171], generative adversarial models and conditional variational autoencoder to generate novel antimicrobial peptides [55,177], and residual networks that predict the tertiary structure of a sequence combined with Monte Carlo sampling to design *de novo* protein structures [178]. These rather *ad-hoc* modeling approaches are now slowly formalized and generalized into mathematical design frameworks [179–181].

Research groups have also started to tackle the long-standing problem of generating sequences that fold into a given 3D structure with deep learning. Early approaches used stacked autoencoder architectures with feed-forward neural networks and a range of local and global structural features [182], while others used convolutional [183] or graph neural networks [184] with self-attention [185] and combined them with a novel neural architecture that can encode geometric features of the problem [186].

A next crucial step is a thorough experimental evaluation of design approaches for a well-defined target to identify strengths and weaknesses in common modeling practices and assumptions. Secondly, many of the current models ignore existing biophysical knowledge. Purely data-driven models seem to have hit an upper ceiling in terms of what can be extracted from protein sequences alone, evident by the minimal improvement in recent large-scale modeling attempts that trained on 2.1 billion sequences [169] over similar models trained on smaller datasets. New approaches that integrate biophysical knowledge through, e.g. physics-informed neural networks [187] or normalizing flows that encode intrinsic invariances of protein structures [188] might further improve such deep generative models by restricting their learned latent space to biophysically plausible regions.

In summary, AI has considerably advanced our ability to predict and learn from the structure and fold of proteins, creating a strong foundation for engineering novel therapeutics and generating new proteins with pre-specified functionality entirely from scratch. We foresee that such approaches will further accelerate personalized medicine enabling the generation of therapeutics tailored to individual genomes.

### 3.3. Small molecule design

Similar to proteins, small molecule design has been driven by the success of deep generative models (Figure 3(c)). In contrast to proteins, small molecules cannot tap into long evolutionary trajectories to extract fundamental features. Instead, deep generative models need to directly extract the biophysical rules that make up small molecule structures from large diverse datasets. Representation learning allows AI to uncover latent factors that determine the properties of small molecules and how to tweak these factors to generate novel molecules with desired properties [189–192].

Early attempts represented small molecule structures as simplified molecular-input line-entry system (SMILE) strings and applied common models from natural language processing [193]. However, these naive models often generated invalid SMILE strings, which led to the incorporation of context-free grammars into such models to constrain them to

produce valid strings [194]. Others have directly encoded small molecules as graphs and used junction trees to iteratively generate small molecules [195], or used graph convolutional and attention-based neural architectures [196–199]. Graph-based approaches were recently used to identify highly potent novel antimicrobial drugs to treat multi-resistant bacterial strains [56]. Some groups also transformed small molecules into 2D [200] or 3D images [201] and used deconvolutional architectures as generators.

The learned latent representations can be used in conjunction with optimization techniques such as Bayesian optimization [193,202] and reinforcement learning [203,204] to generate novel small molecules with optimized properties. Initial results have relied partially or completely on user-specified scoring functions to guide the optimization, but recent works integrate predictors of molecular properties into the reward signal [205] and extend the design to multi-criterion optimization [206]. By departing from continuous latent embeddings and using chemical domain knowledge, novel reinforcement learning methods were able to create new drugs by modifying existing molecules, without the need for expensive pre-training and massive datasets [189]. Such reinforcement learning-based models were recently successfully used to generate highly potent DDR1 kinases inhibitors [57], demonstrating the power of such AI-driven approaches.

A parallel research direction emerged from conditional generative models, in which additional information about the compound is made available to the model while learning its latent representation. By learning to jointly predict molecule properties and structure from embeddings, the search for drugs with pre-specified properties can be performed more efficiently [190–192].

Despite the increase of new AI-based small molecule design approaches, a common set of validation criteria has not yet emerged. Two recent benchmarking platforms addressed these issues [207,208]. Polykovskiy et al. [208] evaluated generative models based on their capability to produce structures similar to the training set, while Brown et al. [207] probed the models' ability to replicate the physicochemical property distributions of a reference set, and to generate novel molecules that optimize single or multiple criteria jointly.

In summary, small molecule design has greatly benefited from AI's capabilities of learning latent representations that drive the functional properties of such molecules and exploiting the euclidean structure of the resulting embedding manifold to improve said properties. This enabled efficient library design for high-throughput drug screenings which can ultimately translate into a significant increase in the success rate of downstream clinical trials, as poor drug candidates could be reliably identified and discarded in silico. Most generative models are trained end-to-end and only learn about physical plausibility implicitly, thus the produced molecules necessitate post-hoc structural fine-tuning using molecular dynamics. Similar to generative models for proteins, more physics-informed networks that constrain the latent space to regions yielding physicochemical viable solutions might improve deep learning-based *de novo* design. Even though the application of such models in a personalized setting has yet to be shown, we envision that with improved pharmacogenomics screening datasets, models will be soon developed

that can generate novel molecules conditioned on mutational changes of their target protein.

## 4. Expert opinion

Artificial intelligence has enabled precision medicine through advancements in biomarker discovery, drug repurposing, combination- and drug design. The successful deployment of the discussed AI-based methods in drug discovery and drug development will enable more nuanced and iterative processes compared to the conventional linear drug development pipelines that may, ultimately, reduce R&D costs and increase the success rate of clinical trials, thus making the pharmaceutical sector as a whole more efficient. However, this development comes along with many challenges and caveats.

AI-driven approaches in early drug discovery are challenging to validate. Despite the fast increase in the number of computational models for vaccine, protein, and small-molecule design, only a minority have experimental validations, or are validated on a common set of benchmarking samples. This may be overcome with regular crowd-sourcing competitions, such as the Dialog for Reverse Engineering Assessments and Methods (DREAM) challenges [209], enabling an unbiased evaluation of computational methods, and fostering collaborations in a quickly emerging scientific field.

AI-aided drug design is currently mostly used as a pre-processing step to reduce the number of compounds or alterations to experimentally test. Yet, to fully exploit the potential of machine learning models in drug discovery, we envision a AI-driven experimental design and a closed-feedback loop to continuously learn from newly generated results. Microfluidic drug development systems [210,211] and lab robotics will play a key component to reach the high degree of automation and throughput to realize such an AI-driven approach. In order to make an impact on drug discovery, these applications will need to demonstrate their ability to arrive at a potential drug candidate faster and more efficiently than in current pipelines.

One of the reasons behind the success of AI is its ability to automatically uncover complex non-linear relationships from heterogeneous sources of data, thus reaching superior performances in prediction tasks. While this is a pivotal aspect of any AI application, it is not enough to motivate broad use of these techniques in pre-clinical settings for precision medicine. Given the limitations imposed by simplified models such as cancer cell lines to model *in vivo* systems, AI models trained on *in vitro* data will only be predictive in patients, if they are customized to capture features which are directly transferable to human samples.

Moreover, state-of-the-art patient stratification in precision medicine still mostly relies on univariate biomarkers. In order to address this, AI seeks to extend our toolbox of established biomarkers. Both supervised and unsupervised methods offer new avenues to identify complex drug- and disease signatures by recognizing patterns across different layers of molecular, chemical, and clinical information.

Unfortunately, the broad adoption in daily clinical use is still lacking behind. The deployment of data-driven biomarkers as companion diagnostics in precision medicine will certainly require new designs of clinical trials, such as master protocols [212] or adaptive studies that are designed for optimizing the biomarker-drug co-development process [213]. Computational systems that can match patients automatically to specialized trials depending on their genetic makeup [214] can help to quickly find appropriate precision trials and ease the burden to reach appropriate cohort sizes.

Besides that, the slow adoption of AI methods for precision medicine in the clinical development of drugs has many technological and regulatory reasons. For one, technical batch effects and the lack of standardized protocols between many datasets may induce biases in medical data, which, if used to design novel algorithms, might be detrimental for generalizable and highly accurate models. Moreover, the imbalanced composition of individuals within patient data may also contain biases based on demographics. The field of fairness in AI focuses on mathematically identifying and addressing such ingrained biases in data, and has started to attract tremendous attention in recent years [215]. We anticipate that this field will become an integral part in high-stakes application areas of AI in healthcare.

Secondly, the data used for drug development in precision medicine comprises highly personalized information such as genetic and clinical variables of an individual, making data privacy an important issue. Often such data cannot easily be shared, or cannot physically leave a specific location, which makes the development of AI applications difficult as many valuable data sets are inaccessible and protected in data silos. Federated learning tackles this problem [216] by enabling decentralized training and prediction of machine learning models without data ever leaving its physical location. Federated learning will enable a connected healthcare network helping physicians to find similar patient cases and inform their treatment decisions without ever needing to exchange sensitive information. It will also enable new business models in which specialized companies can offer AI-based analytics for healthcare and biomedical analysis without needing direct access to the data. However, its broader adoption is hampered by technical and legal challenges and is highly dependent on homogenized healthcare records across clinical networks. We believe that in time, as this technology matures, these challenges will also be overcome.

AI's stigmatization as a 'black-box' approach is also hampering advances in precision medicine as legal liabilities of such 'black-box' predictions in clinics are not satisfactorily answered yet. AI ought to be interpretable, trustworthy, robust, and transparent [217], as encouraged by the Ethics guidelines for trustworthy AI [218]. To ensure that these criteria are met, a growing body of research has been focusing on the areas of interpretable AI [219] and uncertainty quantification [220], which are crucial issues to solve in order to deploy and integrate computational models in the decision-making process in high-risk scenarios. Efforts in this direction have tried to combine state-of-the-art neural network models with different AI-based solutions, such as

symbolic reasoning algorithms, to increase model interpretability while additionally overcoming other traditional weaknesses of black-box models such as the need for big datasets for training and poor generalization capabilities. For example, concepts stemming from symbolic AI gave rise to hybrid AI models that are able to leverage ontologies for guided training and thus may enable human experts to understand and explain predictions [217]. Such modifications of traditional deep learning methods, together with the definition of robust signatures in conjunction with biomedical expert domain knowledge, will potentially enable biological interpretation of these deep molecular profiles, thus turning the aforementioned 'black-boxes' into 'white-boxes'. Therefore, we envision that integration of expert knowledge into AI-approaches and the application of interpretable AI methods will drive the next generation of precision medicine.

Eventually, personalized medicine, in which treatments are truly tailored to the individual, will be possible in a routine clinical setting through a high degree of automation including sequencing, analysis of the biological data, and patient-specific drug development. Here, AI will play a key role in the development of autonomous processes, which could ultimately and affordably bring personalized treatments into clinical practice. Such integrated analysis and design pipelines are currently being established in larger clinical centers. But this also highlights the issues of personalized approaches: they are more laborious to produce and raise ethical questions regarding availability and coverage of costs as large core facilities are required.

As AI and data collection technologies mature, it will eventually be possible to model a significant portion of biological systems in patients, representing the patient's own 'digital twin', a computer model that can be used to estimate the patient's response to a given therapy. Such AI-supported digital twins will allow doctors to quickly tailor a treatment plan to the individual circumstances of their patients. The availability of such information will impact the pharmaceutical sector and will drive the translation from biomarker discovery to drug design in a fully integrated and automated way, lowering potential costs along the drug development pipeline, and improving the efficiency of clinical trials. This transformation, already under way in precision oncology, will soon hit other medical fields and will lead to a radical transformation of the pharmaceutical sector.

In conclusion, precision medicine, especially personalized medicine, is impossible to realize in clinical practice without the aid of advanced AI methods emerging in drug discovery and development. Despite the advancements in treatment strategies and the increase in genomics and molecular information available, the drug development pipeline is still a slow and inefficient process. Acceleration and adoption of common drug development practices for precision and personalized medicine are one of the great challenges facing medical research and development to date. The shift toward a data-driven healthcare system will have far-reaching implications for patients, clinicians, and the pharmaceutical industry. Many technical and regulatory hurdles have to be overcome to make an AI-driven precision medicine approach a reality. However, we strongly believe that the positive societal impact will be profound.

## Declaration of interest

The author(s) have no other relevant affiliations or financial involvement with any organization or entity with a financial interest in or financial conflict with the subject matter or materials discussed in the manuscript apart from those disclosed.

## Reviewer disclosures

Peer reviewers on this manuscript have no relevant financial or other relationships to disclose.

## ORCID

Michael P. Menden http://orcid.org/0000-0003-0267-5792

## References

Papers of special note have been highlighted as either of interest (•) or of considerable interest (••) to readers.

1. Wouters OJ, McKee M, Luyten J. Estimated research and development investment needed to bring a new medicine to market, 2009-2018. JAMA. 2020;323(9):844–853.
2. Wong CH, Siah KW, Lo AW. Estimation of clinical trial success rates and related parameters. Biostatistics. 2019;20(2):273–286.
3. Garraway LA, Verweij J, Ballman KV. Precision oncology: an overview. J Clin Oncol. 2013;31(15):1803–1805.
4. Strianese O, Rizzo F, Ciccarelli M, et al. Precision and personalized medicine: how genomic approach improves the management of cardiovascular and neurodegenerative disease. Genes (Basel). 2020;11(7):11.
5. Tang WHW, Wilcox JD, Jacob MS, et al. Comprehensive diagnostic evaluation of cardiovascular physiology in patients with pulmonary vascular disease: insights from the PVDOMICS program. Circ Heart Fail. 2020;13(3):e006363. .
6. Angwin C, Jenkinson C, Jones A, et al. TriMaster: randomised double-blind crossover study of a DPP4 inhibitor, SGLT2 inhibitor and thiazolidinedione as second-line or third-line therapy in patients with type 2 diabetes who have suboptimal glycaemic control on metformin treatment with or without a sulfonylurea-a MASTERMIND study protocol. BMJ Open. 2020;10 (12):e042784. .
7. Hampel H, Williams C, Etcheto A, et al. A precision medicine framework using artificial intelligence for the identification and confirmation of genomic biomarkers of response to an Alzheimer's

disease therapy: analysis of the blarcamesine (ANAVEX2-73) Phase 2a clinical study. Alzheimers Dement. 2020;6:e12013.

8. Morello G, Salomone S, D'Agata V, et al. From multi-omics approaches to precision medicine in amyotrophic lateral sclerosis. Front Neurosci. 2020;14:577755.

9. Morello G, Guarnaccia M, Spampinato AG, et al. Integrative multi-omic analysis identifies new drivers and pathways in molecularly distinct subtypes of ALS. Sci Rep. 2011;79(1):9968. .

10. Malone ER, Oliva M, Sabatini PJB, et al. Molecular profiling for precision cancer therapies. Genome Med. 2020;12(1):8.

11. Mutz K-O, Heilkenbrinker A, Lönne M, et al. Transcriptome analysis using next-generation sequencing. Curr Opin Biotechnol. 2013;24(1):22–30.

12. Mensaert K, Denil S, Trooskens G, et al. Next-generation technologies and data analytical approaches for epigenomics. Environ Mol Mutagen. 2014;55(3):155–170.

13. Ang MY, Low TY, Lee PY, et al. Proteogenomics: from next-generation sequencing (NGS) and mass spectrometry-based proteomics to precision medicine. Clin Chim Acta. 2019;498:38–46.

14. Hughes JP, Rees S, Kalindjian SB, et al. Principles of early drug discovery. Br J Pharmacol. 2011;162(6):1239–1249.

15. Dugger SA, Platt A, Goldstein DB. Drug development in the era of precision medicine. Nat Rev Drug Discov. 2018;17(3):183–196.

16. Ginsburg GS, McCarthy JJ. Personalized medicine: revolutionizing drug discovery and patient care. Trends Biotechnol. 2001;19(12):491–496.

17. Guinney J, Dienstmann R, Wang X, et al. The consensus molecular subtypes of colorectal cancer. Nat Med. 2015;21(11): 1350–1356.
 •• **One of the most prominent cancer subtyping classifiers that is currently investigated in a clinical setting.**

18. Biton A, Bernard-Pierrot I, Lou Y, et al. Independent component analysis uncovers the landscape of the bladder tumor transcriptome and reveals insights into luminal and basal subtypes. Cell Rep. 2014;9(4):1235–1245 .

19. Wang Z, Wang Y. Extracting a biologically latent space of lung cancer epigenetics with variational autoencoders. BMC Bioinformatics. 2019;20(S18):568.

20. Zhang L, Lv C, Jin Y, et al. Deep learning-based multi-omics data integration reveals two prognostic subtypes in high-risk neuroblastoma. Front Genet. 2018;9:477.

21. Lee T-Y, Huang K-Y, Chuang C-H, et al. Incorporating deep learning and multi-omics autoencoding for analysis of lung adenocarcinoma prognostication. Comput Biol Chem. 2020;87:107277.

22. Chen R, Yang L, Goodison S, et al. Deep-learning approach to identifying cancer subtypes using high-dimensional genomic data. Bioinformatics. 2020;36(5):1476–1483.

23. Fang C, Xu D, Su J, et al. DeePaN: deep patient graph convolutional network integrating clinico-genomic evidence to stratify lung cancers for immunotherapy. Npj Digital Med. 2021;4(1). 10.1038/s41746-021-00381-z

24. El-Manzalawy Y CCA based multi-view feature selection for multi-omics data integration. 2018 IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB). 2018. pp. 1–8.

25. Argelaguet R, Velten B, Arnol D, et al. Multi-omics factor analysis-a framework for unsupervised integration of multi-omics data sets. Mol Syst Biol. 2018;14(6):e8124. .

26. Zhang X, Zhang J, Sun K, et al. Integrated multi-omics analysis using variational autoencoders: application to pan-cancer classification. 2019 IEEE International Conference on Bioinformatics and Biomedicine (BIBM). 2019. pp. 765–769.

27. Azuaje F. Computational models for predicting drug responses in cancer research. Brief Bioinform. 2017;18(5):820–829.

28. Menden MP, Iorio F, Garnett M, et al. Machine learning prediction of cancer cell sensitivity to drugs based on genomic and chemical properties. PLoS One. 2013;8(4): e61318.
 • **First attempts to use drug-based features for improving in vitro drug response predictions.**

29. Wang L, Li X, Zhang L, et al. Improved anticancer drug response prediction in cell lines using matrix factorization with similarity regularization. BMC Cancer. 2017;17(1):513.

30. Ahmadi Moughari F, Eslahchi C. ADRML: anticancer drug response prediction using manifold learning. Sci Rep. 2020;10(1):14245.

31. Chang Y, Park H, Yang H-J, et al. Cancer drug response profile scan (cdrscan): a deep learning model that predicts drug effectiveness from cancer genomic signature. Sci Rep. 2018;8(1):8857. .

32. Rampášek L, Hidru D, Smirnov P, et al. Dr.VAE: improving drug response prediction via modeling of drug perturbation effects. Bioinformatics. 2019;35(19):3743–3751. .
 •• **First variational autoencoder for drug response prediction incorporating drug gene expression signatures.**

33. Manica M, Oskooei A, Born J, et al. Toward explainable anticancer compound sensitivity prediction via multimodal attention-based convolutional encoders. Mol Pharm. 2019;16(12):4797–4806.
 • **One of the first attention-based convolutional encoders allowing explainable drug response prediction.**

34. Kuenzi BM, Park J, Fong SH, et al. Predicting Drug Response and Synergy Using a Deep Learning Model of Human Cancer Cells. Cancer Cell. 2020;38(5): 672–684.e6.
 •• **First efforts for using interpretable VNNs for in vitro drug response prediction with a focus on patient data.**

35. Marcus G. Deep learning: a critical appraisal. arXiv [cs.AI]; 2018. Available: http://arxiv.org/abs/1801.00631

36. Stathias V, Jermakowicz AM, Maloof ME, et al. Drug and disease signature integration identifies synergistic combinations in glioblastoma. Nat Commun. 2018;9(1):5315. .

37. Menden MP, Wang D, Mason MJ, et al. Community assessment to advance computational prediction of cancer drug combinations in a pharmacogenomic screen. Nat Commun. 2019;10(1): 2674.
 •• **First large scale community-driven benchmarking study for in vitro drug synergy prediction.**

38. Xia F, Shukla M, Brettin T, et al. Predicting tumor cell line response to drug pairs with deep learning. BMC Bioinformatics. 2018;19(S18):486. .

39. Preuer K, Lewis RPI, Hochreiter S, et al. DeepSynergy: predicting anti-cancer drug synergy with Deep Learning. Bioinformatics. 2018;34(9):1538–1546. .
 • **One of the earlier studies for using deep learning for predicting potent drug combinations.**

40. Shoemaker RH. The NCI60 human tumour cell line anticancer drug screen. Nat Rev Cancer. 2006;6(10):813–823.

41. Li J, Zheng S, Chen B, et al. A survey of current trends in computational drug repositioning. Brief Bioinform. 2016;17(1):2–12.

42. Yella JK, Yaddanapudi S, Wang Y, et al. Changing trends in computational drug repositioning. Pharmaceuticals. 2018;11(2):57.

43. Ott PA, Hu-Lieskovan S, Chmielowski B, et al. A Phase Ib trial of personalized neoantigen therapy plus anti-PD-1 in patients with advanced melanoma, non-small cell lung cancer, or bladder cancer. Cell. 2020;183(2):347–362.e24. .

44. Hilf N, Kuttruff-Coqui S, Frenzel K, et al. Actively personalized vaccination trial for newly diagnosed glioblastoma. Nature. 2019;565(7738):240–245. .

45. Keskin DB, Anandappa AJ, Sun J, et al. Neoantigen vaccine generates intratumoral T cell responses in phase Ib glioblastoma trial. Nature. 2019;565(7738):234–239. .

46. Ott PA, Hu Z, Keskin DB, et al. An immunogenic personal neoantigen vaccine for patients with melanoma. Nature. 2017;547(7662): 217–221.
 •• **One of the first personalized cancer vaccine trials using Machine Learning to identify neoepitopes.**

47. Sahin U, Derhovanessian E, Miller M, et al. Personalized RNA mutanome vaccines mobilize poly-specific therapeutic immunity against cancer. Nature. 2017;547(7662): 222–226.
 •• **One of the first personalized cancer vaccine trials using Machine Learning to identify neoepitopes.**

48. Waldman AD, Fritz JM, Lenardo MJ. A guide to cancer immunotherapy: from T cell basic science to clinical practice. Nat Rev Immunol. 2020;20:651–668.

49. Pfeifer N, Kohlbacher O. Multiple instance learning allows MHC Class II epitope predictions across alleles. In: Crandall K.A., Lagergren J, editors. Algorithms in bioinformatics. Springer Berlin Heidelberg; 2008. p. 210–221. Available from: https://doi.org/10.1007/978-3-540-87361-7_18

50. Racle J, Michaux J, Rockinger GA, et al. Robust prediction of HLA class II epitopes by deep motif deconvolution of immunopeptidomes. Nat Biotechnol. 2019;37(11):1283–1286. .

51. Reynisson B, Alvarez B, Paul S, et al. NetMHCpan-4.1 and NetMHCIIpan-4.0: improved predictions of MHC antigen presentation by concurrent motif deconvolution and integration of MS MHC eluted ligand data. Nucleic Acids Res. 2020;48(W1):W449–W454. .
•• The leading MHC binding affinity prediction methods with incorporation of MHC ligandomics data.

52. Schubert B, Kohlbacher O. Designing string-of-beads vaccines with optimal spacers. Genome Med. 2016;8(1):9.

53. Toussaint NC, Dönnes P, Kohlbacher O. A mathematical framework for the selection of an optimal set of peptides for epitope-based vaccines. PLoS Comput Biol. 2008;4(12):e1000246. .
•• The first mathematical model for epitope-based vaccine design that guaranteed an global optimal solution.

54. Vider-Shalit T, Raffaeli S, Louzoun Y. Virus-epitope vaccine design: informatic matching the HLA-I polymorphism to the virus genome. Mol Immunol. 2007;44(6):1253–1261.

55. Müller AT, Hiss JA, Schneider G. Recurrent neural network model for constructive peptide design. J Chem Inf Model. 2018;58(2):472–479.

56. Stokes JM, Yang K, Swanson K, et al. A deep learning approach to antibiotic discovery. Cell. 2020;180(4):688–702.e13. .

57. Zhavoronkov A, Ivanenkov YA, Aliper A, et al. Deep learning enables rapid identification of potent DDR1 kinase inhibitors. Nat Biotechnol. 2019;37(9): 1038–1040.
•• One of the first examples of fully Ai-driven small molecule drug design.

58. Schubert B, Schärfe C, Dönnes P, et al. Population-specific design of de-immunized protein biotherapeutics. PLoS Comput Biol. 2018;14(3):e1005983. .
•• One of the first precision protein re-engineering approaches.

59. Paul D, Sanap G, Shenoy S, et al. Artificial intelligence in drug discovery and development. Drug Discov Today. 2021;26(1):80–93.

60. Mazzocchi F. Could Big Data be the end of theory in science? A few remarks on the epistemology of data-driven science. EMBO Rep. 2015;16(10):1250–1255.

61. Gerke S, Minssen T, Cohen G. Ethical and legal challenges of artificial intelligence-driven healthcare. Artificial Intelligence in Healthcare. 2020;295–336. DOI:10.1016/B978-0-12-818438-7.00012-5

62. Torkzadehmahani R, Nasirigerdeh R, Blumenthal DB, et al. Privacy-preserving Artificial Intelligence Techniques in Biomedicine. arXiv preprint arXiv:2007.11621; 2020. Available http://arxiv.org/abs/2007.11621

63. Yu K-H, Beam AL, Kohane IS. Artificial intelligence in healthcare. Nat Biomed Eng. 2018;2(10):719–731.

64. Litjens G, Kooi T, Bejnordi BE, et al. A survey on deep learning in medical image analysis. Med Image Anal. 2017;42:60–88.

65. Teare P, Fishman M, Benzaquen O, et al. Malignancy detection on mammography using dual deep convolutional neural networks and genetically discovered false color input enhancement. J Digit Imaging. 2017;30(4):499–505.

66. Nemati S, Holder A, Razmi F, et al. An interpretable machine learning model for accurate prediction of sepsis in the ICU. Crit Care Med. 2018;46(4):547–553.

67. Shortliffe EH, Sepúlveda MJ. Clinical decision support in the era of artificial intelligence. JAMA. 2018;320(21):2199–2200.

68. Russell, Stuart J. (Stuart Jonathan). Artificial Intelligence : a Modern Approach. Upper Saddle River, N.J. :Prentice Hall, 2010.

69. Breiman L. Random Forests. Machine Learning. 2001;5–32. DOI:10.1023/A:1010933404324.

70. Cortes C, Vapnik V. Support-vector networks. Mach Learn. 1995;20(3):273–297.

71. Oord, Aaron van den, et al. WaveNet: a generative model for raw audio. 2016. Available http://arxiv.org/abs/1609.03499

72. Krizhevsky A, Sutskever I, Hinton GE. ImageNet classification with deep convolutional neural networks. Commun ACM. 2017;60(6):84–90.

73. Miotto R, Wang F, Wang S, et al. Deep learning for healthcare: review, opportunities and challenges. Brief Bioinform. 2018;19(6):1236–1246.

74. Schmidhuber J. Deep learning in neural networks: an overview. Neural Netw. 2015;61:85–117.

75. Lauritsen SM, Kristensen M, Olsen MV, et al. Explainable artificial intelligence model to predict acute critical illness from electronic health records. Nat Commun. 2020;11(1):3852. .

76. Kato S, Subbiah V, Kurzrock R. Counterpoint: successes in the Pursuit of Precision Medicine: biomarkers Take Credit. J Natl Compr Canc Netw. 2017;15(7):863–866.

77. Menden MP, Casale FP, Stephan J, et al. The germline genetic component of drug sensitivity in cancer cell lines. Nat Commun. 2018;9(1):3385. .

78. Azuaje F. Artificial intelligence for precision oncology: beyond patient stratification. NPJ Precis Oncol. 2019;3(1):6.

79. Zhao L, Lee VHF, Ng MK, et al. Molecular subtyping of cancer: current status and moving toward clinical applications. Brief Bioinform. 2019;20(2):572–584.

80. De Sousa E, Melo F, Vermeulen L, et al. Cancer heterogeneity–a multifaceted view. EMBO Rep. 2013;14(8):686–695.

81. Loree JM, Pereira AAL, Lam M, et al. Classifying colorectal cancer by tumor location rather than sidedness highlights a continuum in mutation profiles and consensus molecular subtypes. Clin Cancer Res. 2018;24(5):1062–1072. .

82. Jaber MI, Song B, Taylor C, et al. A deep learning image-based intrinsic molecular subtype classifier of breast tumors reveals tumor heterogeneity that may affect survival. Breast Cancer Res. 2020;22(1):12. .

83. Mazurowski MA, Grimm LJ, Zhang J, et al. Recurrence-free survival in breast cancer is associated with MRI tumor enhancement dynamics quantified using computer algorithms. Eur J Radiol. 2015;84(11):2117–2122. .

84. Yaeger R, Chatila WK, Lipsyc MD, et al. Clinical sequencing defines the genomic landscape of metastatic colorectal cancer. Cancer Cell. 2018;33(1):125–136.e3. .

85. Mooi JK, Wirapati P, Asher R, et al. The prognostic impact of consensus molecular subtypes (CMS) and its predictive effects for bevacizumab benefit in metastatic colorectal cancer: molecular analysis of the AGITG MAX clinical trial. Ann Oncol. 2018;29(11):2240–2246. .

86. Moffitt RA, Marayati R, Flate EL, et al. Virtual microdissection identifies distinct tumor- and stroma-specific subtypes of pancreatic ductal adenocarcinoma. Nat Genet. 2015;47(10):1168–1178. .

87. Orth M, Metzger P, Gerum S, et al. Pancreatic ductal adenocarcinoma: biological hallmarks, current status, and future perspectives of combined modality treatment approaches. Radiat Oncol. 2019;14:141.

88. Tomczak K, Czerwińska P, Wiznerowicz M. The Cancer Genome Atlas (TCGA): an immeasurable source of knowledge. Contemp Oncol. 2015;19:A68–77.

89. Hudson TJ, Anderson W, Artez A, et al.; International Cancer Genome Consortium. International network of cancer genome projects. Nature. 2010;464:993–998.

90. Gerstung M, Jolly C, Leshchiner I, et al. The evolutionary history of 2,658 cancers. Nature. 2020;578(7793):122–128. .

91. Amin SB, Yip W-K, Minvielle S, et al. Gene expression profile alone is inadequate in predicting complete response in multiple myeloma. Leukemia. 2014;28(11):2229–2234. .

92. Argelaguet R, Arnol D, Bredikhin D, et al. MOFA+: a statistical framework for comprehensive integration of multi-modal single-cell data. Genome Biol. 2020;21(1):111. .

93. Yang KD, Multi-Domain UC. Translation by learning uncoupled autoencoders. arXiv [cs.LG]; 2019. Available: http://arxiv.org/abs/1902.03515

94. Iorio F, Knijnenburg TA, Vis DJ, et al. A Landscape of Pharmacogenomic Interactions in Cancer. Cell. 2016;166(3): 740–754.\
•• One of the largest in vitro drug screening efforts in cancer cell lines.

95. Ghandi M, Huang FW, Jané-Valbuena J, et al. Next-generation characterization of the cancer cell line encyclopedia. Nature. 2019;569(7757): 503–508.
•• One of the largest in vitro drug screening efforts in cancer cell lines.

96. Drost J, Clevers H. Organoids in cancer research. Nat Rev Cancer. 2018;18(7):407–418.

97. Gao H, Korn JM, Ferretti S, et al. High-throughput screening using patient-derived tumor xenografts to predict clinical trial drug response. Nat Med. 2015;21(11): 1318–1325.
• One of the largest in vivo drug screening efforts in PDX models.

98. Yu K, Chen B, Aran D, et al. Comprehensive transcriptomic analysis of cell lines as models of primary tumors across 22 tumor types. Nat Commun. 2019;10(1):3574. .

99. Najgebauer H, Yang M, Francies HE, et al. CELLector: genomics-Guided Selection of Cancer In Vitro Models. Cell Syst. 2020;10(5):424–432.e6. .

100. Vamathevan J, Clark D, Czodrowski P, et al. Applications of machine learning in drug discovery and development. Nat Rev Drug Discov. 2019;18(6):463–477. .

101. Tsherniak A, Vazquez F, Montgomery PG, et al. Defining a cancer dependency map. Cell. 2017;170(3):564–576.e16. .

102. Behan FM, Iorio F, Picco G, et al. Prioritization of cancer therapeutic targets using CRISPR–Cas9 screens. Nature. 2019;568(7753): 511–516.
•• First and biggest CRISPR-Cas9 screen to this day, enabling data-driven discovery of novel therapeutic targets.

103. Costello JC, Heiser LM, Georgii E, et al. A community effort to assess and improve drug sensitivity prediction algorithms. Nat Biotechnol. 2014;32(12): 1202–1212.
•• The earliest community-driven benchmarking study for machine learning drug response prediction algorithms.

104. Turki T, Wei Z, Wang JTL. Transfer learning approaches to improve drug sensitivity prediction in multiple myeloma patients. IEEE Access. 2017;5:7381–7393.

105. Subramanian A, Narayan R, Corsello SM, et al. A next generation connectivity map: L1000 platform and the first 1,000,000 profiles. Cell. 2017;171(6): 1437–1452.e17.
• The L1000 platform is the largest drug perturbation screen until this day, enabling many drug repurposing studies.

106. Zhao W, Li J, Chen M-JM, et al. Large-scale characterization of drug responses of clinically relevant proteins in cancer cell lines. Cancer Cell. 2020;38(6):829–843.e4.

107. Musa A, Ghoraie LS, Zhang S-D, et al. A review of connectivity map and computational approaches in pharmacogenomics. Brief Bioinform. 2018;19(3):506–523. .

108. Gaulton A, Bellis LJ, Bento AP, et al. ChEMBL: a large-scale bioactivity database for drug discovery. Nucleic Acids Res. 2012;40(D1): D1100–7. .

109. Kim S, Thiessen PA, Bolton EE, et al. PubChem substance and compound databases. Nucleic Acids Res. 2016;44(D1):D1202–13. .

110. Cheng F, Liu C, Jiang J, et al. Prediction of drug-target interactions and drug repositioning via network-based inference. PLoS Comput Biol. 2012;8(5):e1002503. .

111. Emig D, Ivliev A, Pustovalova O, et al. Drug target prediction and repositioning using an integrated network-based approach. PLoS One. 2013;8(4):e60618. .

112. Mayr A, Klambauer G, Unterthiner T, et al. Large-scale comparison of machine learning methods for drug target prediction on ChEMBL. Chem Sci. 2018;9(24):5441–5451. .

113. Peyvandipour A, Saberian N, Shafi A, et al. A novel computational approach for drug repurposing using systems biology. Bioinformatics. 2018;34(16):2817–2825.

114. Oh M, Ahn J, Yoon Y. A network-based classification model for deriving novel drug-disease associations and assessing their molecular actions. PLoS One. 2014;9(10):e111668.

115. Brown AS, Kong SW, Kohane IS, et al. ksRepo: a generalized platform for computational drug repositioning. BMC Bioinformatics. 2016;17(1):78.

116. Jadamba E, Shin M. A Systematic Framework for Drug Repositioning from Integrated Omics and Drug Phenotype Profiles Using Pathway-Drug Network. Biomed Res Int. 2016;2016:7147039.

117. Luo Y, Zhao X, Zhou J, et al. A network integration approach for drug-target interaction prediction and computational drug repositioning from heterogeneous information. Nat Commun. 2017;8 (1):573. .

118. Aliper A, Plis S, Artemov A, et al. Deep learning applications for predicting pharmacological properties of drugs and drug repurposing using transcriptomic data. Mol Pharm. 2016;13(7):2524–2530.

119. Wan F, Hong L, Xiao A, et al. NeoDTI: neural integration of neighbor information from a heterogeneous network for discovering new drug–target interactions. Bioinformatics. 2018;35 (1):104–111.

120. Li YY, Jones SJ. Drug repositioning for personalized medicine. Genome Med. 2012;4(3):27.

121. Gysi DM, Do Valle Í, Zitnik M, et al. Network medicine framework for identifying drug repurposing opportunities for COVID-19. ArXiv; 2020. Available: https://www.ncbi.nlm.nih.gov/pubmed/32550253.
• New computational network-based repurposing strategy for the current SARS-CoV-2 pandemic.

122. Wang X, Zhang H, Chen X. Drug resistance and combating drug resistance in cancer. Cancer Drug Resist. 2019;2:141–160.

123. Ayestaran I, Galhoz A, Spiegel E, et al. Identification of Intrinsic Drug Resistance and Its Biomarkers in High-Throughput Pharmacogenomic and CRISPR Screens. Patterns Prejudice. 2020;1 (5):100065.

124. Liu H, Zhang W, Zou B, et al. DrugCombDB: a comprehensive database of drug combinations toward the discovery of combinatorial therapy. Nucleic Acids Res. 2020;48(D1):D871–D881.

125. Eroglu Z, Ribas A. Combination therapy with BRAF and MEK inhibitors for melanoma: latest evidence and place in therapy. Ther Adv Med Oncol. 2016;8(1):48–56.

126. Qian T, Zhu S, Hoshida Y. Use of big data in drug development for precision medicine: an update. Expert Rev Precis Med Drug Dev. 2019;4(3):189–200.

127. Bai R, Lv Z, Xu D, et al. Predictive biomarkers for cancer immunotherapy with immune checkpoint inhibitors. Biomark Res. 2020;8 (1):34.

128. Koury J, Lucero M, Cato C, et al. Immunotherapies: exploiting the Immune System for Cancer Treatment. J Immunol Res. 2018;2018:9585614.

129. Ren Y, Chen X, Feng M, et al. Gaussian process: a promising approach for the modeling and prediction of peptide binding affinity to MHC proteins. Protein Pept Lett. 2011;18(7):670–678.

130. Bassani-Sternberg M, Chong C, Guillaume P, et al. Deciphering HLA-I motifs across HLA peptidomes improves neo-antigen predictions and identifies allostery regulating HLA specificity. PLoS Comput Biol. 2017;13(8): e1005725.
• One of the first methods that incorporated MHC ligandomics data with automatic MHC deconvolution.

131. O'Donnell TJ, Rubinsteyn A, Laserson U. MHCflurry 2.0: improved Pan-Allele Prediction of MHC Class I-Presented Peptides by Incorporating Antigen Processing. Cell Syst. 2020;11(1):42–48.e7.

132. Shao XM, Bhattacharya R, Huang J, et al. High-throughput prediction of MHC class I and II neoantigens with MHCnuggets. Cancer Immunol Res. 2020;8(3):396–408. .

133. Freudenmann LK, Marcu A, Stevanović S. Mapping the tumour human leukocyte antigen (HLA) ligandome by mass spectrometry. Immunology. 2018;154(3):331–345.

134. Vita R, Mahajan S, Overton JA, et al. The Immune Epitope Database (IEDB): 2018 update. Nucleic Acids Res. 2019;47(D1):D339–D343. .

135. Marcu A, Bichmann L, Kuchenbecker L, et al. The HLA Ligand Atlas: a novel immuno-oncology resource for T-cell antigen discovery. J Clin Orthod. 2020;38:3128.

136. Chen B, Khodadoust MS, Olsson N, et al. Predicting HLA class II antigen presentation through integrated deep learning. Nat Biotechnol. 2019;37(11): 1332–1343.
   • MHC neoepitope prediction method that also includes gene expression information.

137. Hundal J, Kiwala S, McMichael J, et al. pVACtools: a computational toolkit to identify and visualize cancer neoantigens. Cancer Immunol Res. 2020;8(3):409–420. .

138. Schubert B, Walzer M, Brachvogel H-P, et al. FRED 2: an immunoinformatics framework for Python. Bioinformatics. 2016;32 (13):2044–2046.

139. Bjerregaard A-M, Nielsen M, Hadrup SR, et al. MuPeXI: prediction of neo-epitopes from tumor sequencing data. Cancer Immunol Immunother. 2017;66(9):1123–1130.

140. Wells DK, Van Buuren MM, Dang KK, et al. Key parameters of tumor epitope immunogenicity revealed through a consortium approach improve neoantigen prediction. Cell. 2020;183(3):818–834.e13. .
   • Major benchmark to identify properties of neoepitopes that are predictive for their immunogenicity.

141. Glanville J, Huang H, Nau A, et al. Identifying specificity groups in the T cell receptor repertoire. Nature. 2017;547(7661): 94–98.
   •• First two models predicting TCR-specificity based on sequence similarity.

142. Dash P, Fiore-Gartland AJ, Hertz T, et al. Quantifiable predictive features define epitope-specific T cell receptor repertoires. Nature. 2017;547(7661): 89–93.
   •• First two models predicting TCR-specificity based on sequence similarity.

143. Jokinen E, Huuhtanen J, Mustjoki S, et al. Determining epitope specificity of T cell receptors with TCRGP. bioRxiv; 2019. Available: https://www.biorxiv.org/content/10.1101/542332v2.abstract

144. Gielis S, Moris P, Bittremieux W, et al. Detection of Enriched T Cell Epitope Specificity in Full T Cell Receptor Sequence Repertoires. Front Immunol. 2019;10:2820.

145. Jurtz VI, Jessen LE, Bentzen AK, et al. NetTCR: sequence-based prediction of TCR binding to peptide-MHC complexes using convolutional neural networks. 2018. 433706. DOI: 10.1101/433706

146. Fischer DS, Wu Y, Schubert B, et al. Predicting antigen specificity of single T cells based on TCR CDR3 regions. Mol Syst Biol. 2020;16(8): e9416.

147. Schumacher TN, Schreiber RD. Neoantigens in cancer immunotherapy. Science. 2015;348(6230):69–74.

148. Lundegaard C, Buggert M, Karlsson AC, et al. PopCover: a method for selecting of peptides with optimal population and pathogen coverage. Proceedings of the First ACM International Conference on Bioinformatics and Computational Biology. New York, NY, USA: Association for Computing Machinery; 2010. pp. 658–659.

149. Toussaint NC, Maman Y, Kohlbacher O, et al. Universal peptide vaccines–optimal peptide vaccine design based on viral sequence conservation. Vaccine. 2011;29(47):8745–8753.

150. Ali M, Foldvari Z, Giannakopoulou E, et al. Induction of neoantigen-reactive T cells from healthy donors. Nat Protoc. 2019;14(6):1926–1943. .

151. Dorigatti E, Schubert B. Joint epitope selection and spacer design for string-of-beads vaccines. Bioinformatics. 2020;36 (Supplement_2):i643–i650.

152. Dorigatti E, Schubert B, Kouyos RD. Graph-theoretical formulation of the generalized epitope-based vaccine design problem. PLoS Comput Biol. 2020;16(10):e1008237. .
   • First unified mathematical Framework combining all design steps and techniques.

153. Barouch DH, O'Brien KL, Simmons NL, et al. Mosaic HIV-1 vaccines expand the breadth and depth of cellular immune responses in rhesus monkeys. Nat Med. 2010;16(3):319–323. .

154. Wang W, Li R, Deng Y, et al. Protective efficacy of the conserved NP, PB1, and M1 proteins as immunogens in DNA- and vaccinia virus-based universal Influenza A virus vaccines in mice. Clin Vaccine Immunol. 2015;22(6):618–630. .

155. Audran R, Cachat M, Lurati F, et al. Phase I malaria vaccine trial with a long synthetic peptide derived from the merozoite surface protein 3 antigen. Infect Immun. 2005;73(12):8017–8026. .

156. Von Delft A, Donnison TA, Lourenço J, et al. The generation of a simian adenoviral vectored HCV vaccine encoding genetically conserved gene segments to target multiple HCV genotypes. Vaccine. 2018;36(2):313–321. .

157. Yang Z, Bogdan P, Nazarian S. An in silico deep learning approach to multi-epitope vaccine design: a SARS-CoV-2 case study. Sci Rep. 2021;11(1):3238.

158. Marks DS, Colwell LJ, Sheridan R, et al. Protein 3D structure computed from evolutionary sequence variation. PLoS One. 2011;6(12): e28766.
   • One of the earliest papers on evolutionary based protein sequence modeling.

159. Balakrishnan S, Kamisetty H, Carbonell JG, et al. Learning generative models for protein fold families. Proteins. 2011;79(4):1061–1078.

160. Hopf TA, Schärfe CPI, Rodrigues JPGLM, et al. Sequence co-evolution gives 3D contacts and structures of protein complexes. Elife. 2014;3:3.

161. Weigt M, White RA, Szurmant H, et al. Identification of direct residue contacts in protein–protein interaction by message passing. Proc Natl Acad Sci U S A. 2009;106(1):67–72.

162. Green AG, Elhabashy H, Brock KP, et al. Proteome-scale discovery of protein interactions with residue-level resolution using sequence coevolution. bioRxiv; 2021;12(1):1–12. Available: https://www.biorxiv.org/content/10.1101/791293v2.abstract

163. Hopf TA, Ingraham JB, Poelwijk FJ, et al. Mutation effects predicted from sequence co-variation. Nat Biotechnol. 2017;35 (2):128–135. .

164. Riesselman AJ, Ingraham JB, Marks DS. Deep generative models of genetic variation capture the effects of mutations. Nat Methods. 2018;15(10):816–822. .
   • One of the first deep learning based generative models of protein sequences.

165. Hawkins-Hooker A, Depardieu F, Baur S, et al. Generating functional protein variants with variational autoencoders. Cold Spring Harbor Laboratory; 2021;17(2), e1008736. DOI:10.1101/2020.04.07.029264.

166. Amimeur T, Shaver JM, Ketchem RR, et al. Designing feature-controlled humanoid antibody discovery libraries using generative adversarial networks. Cold Spring Harbor Laboratory;2020. DOI:10.1101/2020.04.12.024844.

167. Alley EC, Khimulya G, Biswas S, et al. Unified rational protein engineering with sequence-based deep representation learning. Nat Methods. 2019;16(12):1315–1322.

168. Heinzinger M, Elnaggar A, Wang Y, et al. Modeling aspects of the language of life through transfer-learning protein sequences. BMC Bioinformatics. 2019;20(1):723. .

169. Elnaggar A, Heinzinger M, Dallago C, et al. ProtTrans: towards Cracking the Language of Life's Code Through Self-Supervised Deep Learning and High Performance Computing. arXiv [cs.LG]; 2020. Available http://arxiv.org/abs/2007.06225

170. Rives A, Goyal S, Meier J, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. bioRxiv; 2021;118(15) e2016239118. Available: https://www.biorxiv.org/content/10.1101/622803v1.abstract

171. Riesselman AJ, Shin JE, Kollasch AW, et al. Accelerating protein design using autoregressive generative models. bioRxiv; 2019. Available: https://www.biorxiv.org/content/10.1101/757252v1.abstract

172. Linder J, Bogard N, Rosenberg AB, et al. Network for maximizing fitness and diversity of synthetic DNA and protein sequences. Cell Syst. 2020;11(1):49–62.e16.

173. Senior AW, Evans R, Jumper J, et al. Improved protein structure prediction using potentials from deep learning. Nature. 2020;577 (7792): 706–710.
    •• **A breakthrough using deep learning for protein structure prediction.**

174. Xu J, Mcpartlon M, Li J. Improved protein structure prediction by deep learning irrespective of co-evolution information. Cold Spring Harbor Laboratory; 2020. DOI: 10.1101/2020.10.12.336859.

175. AlQuraishi M. End-to-End differentiable learning of protein structure. Cell Syst. 2019;8(4):292–301.e3. .
    • **One of the first end-to-end differentiable protein folding approaches.**

176. Ingraham J, Riesselman AJ, Sander C, et al. Learning protein structure with a differentiable simulator. New Orleans, USA: ICLR, 2019.
    • **One of the first end-to-end differentiable protein folding approaches.**

177. Das P, Wadhawan K, Chang O, et al. PepCVAE: semi-Supervised Targeted Design of Antimicrobial Peptide Sequences. arXiv [q-bio. QM]; 2018. Available http://arxiv.org/abs/1810.07743

178. Anishchenko I, Chidyausiku TM, Ovchinnikov S, et al. De novo protein design by deep network hallucination. 2020. DOI:10.1101/2020.07.22.211482.
    • **One of the first applications of fully AI driven de novo protein design using co-evolutionary information.**

179. Brookes DH, Park H, Listgarten J. Conditioning by adaptive sampling for robust design. arXiv [cs.LG]; 2020. Available: http://arxiv.org/abs/1901.10060

180. Kumar A, Levine S. Model Inversion Networks for Model-Based Optimization. arXiv [cs.LG]; 2019. Available: http://arxiv.org/abs/1912.13464

181. Fannjiang C, Listgarten J. Autofocused oracles for model-based design. arXiv [cs.LG]; 2020. Available: http://arxiv.org/abs/2006.08052

182. O'Connell J, Li Z, Hanson J, et al. SPIN2: predicting sequence profiles from protein structures using deep neural networks. Proteins. 2018;86(6):629–633. .

183. Zhang Y, Chen Y, Wang C, et al. ProDCoNN: protein design using a convolutional neural network. Proteins. 2020;88(7):819–829. .

184. Strokach A, Becerra D, Corbi-Verge C, et al. Fast and flexible design of novel proteins using graph neural networks. Cold Spring Harbor Laboratory; 2020;11(4):402–411.e4. DOI:10.1101/868935.

185. Ingraham J, Garg V, Barzilay R, et al. Generative models for graph-based protein design Adv. Neural Inf. Process. Syst. 2019, pp.15820–15831.

186. Jing B, Eismann S, Suriana P, et al. Learning from protein structure with geometric vector perceptrons. arXiv [q-bio.BM]; 2020. Available: http://arxiv.org/abs/2009.01411

187. Raissi M, Perdikaris P, Karniadakis GE. Physics-informed neural networks: a deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. J Comput Phys. 2019;378:686–707.

188. Köhler J, Klein L, Noé F. Equivariant Flows: exact likelihood generative learning for symmetric densities. arXiv preprint arXiv, 2020 2006 02425. Available. http://arxiv.org/abs/2006.02425.

189. Zhou Z, Kearnes S, Li L, et al. Optimization of Molecules via Deep Reinforcement Learning. Sci Rep. 2019;9(1):10752.

190. Li Y, Zhang L, Liu Z. Multi-objective de novo drug design with conditional graph generative model. J Cheminform. 2018;10(1):33.

191. Lim J, Ryu S, Kim JW, et al. Molecular generative model based on conditional variational autoencoder for de novo molecular design. J Cheminform. 2018;10(1):31.

192. Méndez-Lucio O, Baillif B, Clevert D-A, et al. De novo generation of hit-like molecules from gene expression signatures using artificial intelligence. Nat Commun. 2020;11(1):10. .
    • **One of the first frameworks that can generate hit-like small molecules.**

193. Gómez-Bombarelli R, Wei JN, Duvenaud D, et al. Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules. ACS Cent Sci. 2018;4(2): 268–276.
    •• **One of the first end-to-end deep learning methods for optimizing drug formulations.**

194. Kusner MJ, Paige B, Hernández-Lobato JM Grammar Variational Autoencoder. Proceedings of the 34th International Conference on Machine Learning - Volume 70. Sydney, NSW, Australia: JMLR. org; 2017. pp. 1945–1954.
    • **One of the first methods deriving continuous latent representations of molecules.**

195. Jin W, Barzilay R, Jaakkola T. Junction tree variational autoencoder for molecular graph generation. arXiv [cs.LG]; 2018. Available http://arxiv.org/abs/1802.04364.
    • **One of the first graph-based deep learning models of small molecules.**

196. Liao R, Li Y, Song Y, et al. Efficient graph generation with graph recurrent attention networks. Wallach H, Larochelle H, Beygelzimer A, et al., editors. Advances in neural information processing systems 32. Curran Associates, Inc.; 2019. 4255–4265.

197. You J, Liu B, Ying Z, et al. Graph convolutional policy network for goal-directed molecular graph generation. In: Bengio S, Wallach H, Larochelle H, et al., editors. Advances in neural information processing systems 31. Curran Associates, Inc.; 2018. p. 6410–6421.

198. Samanta B, De A, Jana G, et al. Nevae: a deep generative model for molecular graphs. J Mach Learn Res. 2020; Available: https://www.jmlr.org/papers/volume21/19-671/19-671.pdf

199. De Cao N, MolGAN: KT. An implicit generative model for small molecular graphs. arXiv [stat.ML]; 2018. Available: http://arxiv.org/abs/1805.11973

200. Goh GB, Siegel C, Vishnu A, et al. Chemception: a deep neural network with minimal chemistry knowledge matches the performance of expert-developed QSAR/QSPR models. arXiv [stat.ML]; 2017. Available: http://arxiv.org/abs/1706.06689

201. Kuzminykh D, Polykovskiy D, Kadurin A, et al. 3D molecular representations based on the wave transform for convolutional neural networks. Mol Pharm. 2018;15(10):4378–4385. .

202. Griffiths -R-R, Hernández-Lobato JM. Constrained Bayesian optimization for automatic chemical design using variational autoencoders. Chem Sci. 2020;11(2):577–586.

203. Olivecrona M, Blaschke T, Engkvist O, et al. Molecular de-novo design through deep reinforcement learning. J Cheminform. 2017;9(1):48.

204. Guimaraes GL, Sanchez-Lengeling B, Outeiral C, et al. Objective-reinforced generative adversarial networks (ORGAN) for sequence generation models. arXiv [stat.ML]; 2017. Available: http://arxiv.org/abs/1705.10843

205. Popova M, Isayev O, Tropsha A. Deep reinforcement learning for de novo drug design. Sci Adv. 2018;4(7):7885.

206. Ståhl N, Falkman G, Karlsson A, et al. Deep Reinforcement Learning for Multiparameter Optimization in de novo Drug Design. J Chem Inf Model. 2019;59(7):3166–3176.

207. Brown N, Fiscato M, Segler MHS, et al. GuacaMol: benchmarking Models for de Novo Molecular Design. J Chem Inf Model. 2019;59 (3):1096–1108.

208. Polykovskiy D, Zhebrak A, Sanchez-Lengeling B, et al. Molecular Sets (MOSES): a benchmarking platform for molecular generation models. arXiv [cs.LG]; 2018. Available http://arxiv.org/abs/1811.12823

209. Stolovitzky G, Monroe D, Califano A. Dialogue on reverse-engineering assessment and methods: the DREAM of high-throughput pathway inference. Ann N Y Acad Sci. 2007;1115 (1):1–22.

210. Regnault C, Dheeman DS, Hochstetter A. Microfluidic Devices for Drug Assays. High Throughput. 2018;7(2):7.

211. Weng L, Spoonamore JE. Droplet microfluidics-enabled high-throughput screening for protein engineering. Micromachines (Basel). 2019;10(11):10.

212. Woodcock J, LaVange LM, Drazen JM. Master protocols to study multiple therapies, multiple diseases, or both. N Engl J Med. 2017;377(1):62–70.

213. Garralda E, Dienstmann R, Piris-Giménez A, et al. New clinical trial designs in the era of precision medicine. Mol Oncol. 2019;13 (3):549–557.

214. Lindsay J, Del Vecchio Fitz C, Zwiesler Z, et al. MatchMiner: an open source computational platform for real-time matching of cancer patients to precision medicine clinical trials using genomic and clinical criteria. Cold Spring Harbor Laboratory; 2017. 199489. DOI:10.1101/199489.

215. Caton S, Haas C. Fairness in machine learning: a survey. arXiv [cs. LG]; 2020. Available: http://arxiv.org/abs/2010.04053

216. Rieke N, Hancox J, Li W, et al. The future of digital health with federated learning. NPJ Digit Med. 2020;3(1):119. .

217. Kelly CJ, Karthikesalingam A, Suleyman M, et al. Key challenges for delivering clinical impact with artificial intelligence. BMC Med. 2019;17(1):195.

218. Floridi L. Establishing the rules for building trustworthy AI. Nature Mach Intell. 2019;1(6):261–262.

219. Molnar C, Casalicchio G, Bischl B. Interpretable machine learning – a brief history, state-of-the-art and challenges. ECML PKDD 2020 Workshops; 2020. 417–431. DOI:10.1007/978-3-030-65965-3_28

220. Abdar M, Pourpanah F, Hussain S, et al. A review of uncertainty quantification in deep learning: techniques, applications and challenges. arXiv [cs.LG]; 2020. Available http://arxiv.org/abs/2011.06225