

# 1 **AI Approaches for the Discovery and** 2 **Validation of Drug Targets**

3  
4 A. Wenteler<sup>1, 2, 3 \*</sup>, C.P. Cabrera<sup>1, 2, 4</sup>, W. Wei<sup>3</sup>, V. Neduva<sup>3</sup>, M.R. Barnes<sup>1, 2, 4, 5</sup>

5  
6 <sup>1</sup> Digital Environment Research Institute, Queen Mary University of London, London, United Kingdom

7 <sup>2</sup> Centre for Translational Bioinformatics, William Harvey Research Institute, Queen Mary University of London,  
8 London, United Kingdom

9 <sup>3</sup> MSD Discovery Centre, London, United Kingdom

10 <sup>4</sup> NIHR Barts Cardiovascular Biomedical Research Centre, Barts and The London School of Medicine and Dentistry,  
11 Queen Mary University of London, London, United Kingdom

12 <sup>5</sup> The Alan Turing Institute, London, United Kingdom

13  
14 \* Corresponding author: a.wenteler@qmul.ac.uk

## 15 **Abstract**

16 Artificial intelligence (AI) holds immense promise for accelerating and improving all  
17 aspects of drug discovery, not least target discovery and validation. By integrating a  
18 diverse range of biological data modalities, AI enables the accurate prediction of drug  
19 target properties, ultimately illuminating biological mechanisms of disease and guiding  
20 drug discovery strategies. Despite the indisputable potential of AI in drug target  
21 discovery, there are many challenges and obstacles yet to be overcome, including  
22 dealing with data biases, model interpretability and generalisability, and the validation  
23 of predicted drug targets to name a few. By exploring recent advancements in AI, this  
24 review showcases current applications of AI for drug target discovery and offers  
25 perspectives on the future of AI for the discovery and validation of drug targets, paving  
26 the way for the generation of novel and safer pharmaceuticals.

27 **Keywords:** Drug discovery, Drug targets, Artificial intelligence, Machine learning,  
28 Multiomics

29  
30 This peer-reviewed article has been accepted for publication but not yet copyedited or typeset, and so may be subject to change during the production process. The article is considered published and may be cited using its DOI.  
10.1017/pcm.2024.4

This is an Open Access article, distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivatives licence (<http://creativecommons.org/licenses/by-nc-nd/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is unaltered and is properly cited. The written permission of Cambridge University Press must be obtained for commercial re-use or in order to create a derivative work.

## 31 **Background**

32 Historically, drug target discovery and validation has been a laborious and somewhat  
33 haphazard process, heavily reliant on industry standard laboratory models and  
34 analysis procedures (Drews 2000; Huang *et al.* 2004; Materi and Wishart 2007). Most  
35 drug discovery to date has taken a phenotype-first approach focusing on evaluation  
36 of the therapeutic potential of compounds in phenotypic assays, without necessarily  
37 knowing the exact target or mechanism of action (Moffat *et al.* 2017). This approach  
38 relies largely on serendipity, where complex compound libraries, including  
39 phytochemicals, biochemicals and other organic chemistry, are identified for  
40 therapeutic use by chance. Naturally, pharma companies initially sought to improve  
41 their odds by increasing the size and complexity of their compound libraries, and by  
42 the mid-2000s most major pharmaceutical companies had compound libraries in the  
43 range of 1-2 million small molecule entities (SMEs) (Hann and Oprea 2004). However,  
44 the unsustainability of this chemistry arms race has spurred a shift towards a target-  
45 first strategy, which signified a pivotal moment in pharmacological research,  
46 emphasising the importance of thorough understanding and validation of a biological  
47 target before initiation of the drug design process. This paradigm shift marked a  
48 transition from empirical, trial-and-error methods to a more rational and systematic  
49 approach, greatly enhancing the efficiency and effectiveness of drug discovery.  
50 Ironically, although the target-first approach was designed to reduce the complexity  
51 of drug discovery, it may have had the opposite effect, simply highlighting the  
52 challenges of true target validation, leading to over a decade of increased failure in  
53 drug discovery stemming from poorly validated targets (Paul *et al.* 2010; Scannell *et*  
54 *al.* 2012). With an increasing repertoire of biomolecular assays to probe mechanism,  
55 such as CRISPR-Cas9, so-called target deconvolution studies have been conducted.  
56 These studies connect phenotypic to target-first approaches by attempting to elucidate  
57 the mechanism of action of the target upon which a drug acted retrospectively. This  
58 strategy enriches the phenotype-centric drug discovery paradigm with mechanistic  
59 understanding of the observed therapeutic effect and set the groundwork for  
60 integration of phenotype-first and target-first approaches (Terstappen *et al.* 2007).

61  
62 In this review, we define drug targets as biomolecules—primarily proteins, but also  
63 DNA, RNA, or other biomolecular species—that a therapeutic compound can bind to  
64 or modulate. The pool of existing drug targets is limited, and assessments of the  
65 druggable genome, which refers to those genes susceptible to modulation by small

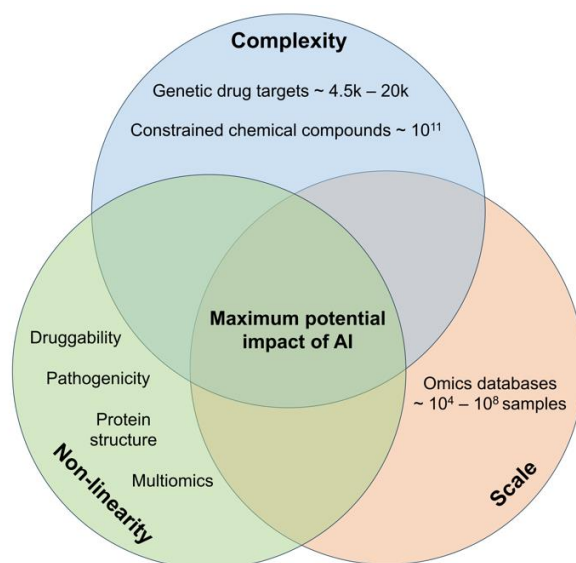
66 molecules, fluctuate. The latest estimate places this number at 4,479 potential  
67 targets, accounting for approximately 22% of protein-coding genes (Finan *et al.*  
68 2017). According to records of the Human Protein Atlas (HPA), there are  
69 approximately 863 FDA approved drug targets (Paananen and Fortino 2020), over  
70 50% of these targets are represented by just four protein families - ion channels,  
71 nuclear receptors, kinases, and G-protein coupled receptors (Bakheet and Doig 2009;  
72 Santos *et al.* 2017). When it comes to finding novel, efficacious, and safer drug  
73 targets, as a general guideline, targets should have a role in disease, limited role in  
74 normal physiology, particularly in critical tissues such as the heart, and ideally should  
75 be druggable with small molecules, although biologic drugs and gene targeted  
76 therapies make almost all targets therapeutically tractable. Furthermore, while a  
77 laboratory resolved 3D protein structure was a prior requirement for drug design, with  
78 the advent of protein structure prediction models, further accelerated by AI  
79 approaches (Baek *et al.* 2021; Jumper *et al.* 2021; Lin *et al.* 2023b), high quality 3D  
80 structures of a wide range of potential drug targets are generally available. This  
81 enables a broader application of *in silico* structure-based drug design. Another  
82 desirable property for a drug target is having multiple binding pockets. By having  
83 multiple potential binding pockets, different conformations of the protein in various  
84 functional states can be targeted. It also provides opportunities for identifying  
85 allosteric inhibitors rather than only targeting the active site. Allosteric sites may offer  
86 better selectivity and provide safety benefits (Abdel-Magid 2015; May *et al.* 2007;  
87 Wagner *et al.* 2016b). Lastly, by understanding the associated pathways of the target,  
88 we gain insight into the processes the target is involved in and thus, what other  
89 biological processes could potentially be affected. This can help the assessment of  
90 potential off-target effects.

91  
92 Despite the great progress in drug discovery, the process is still burdened by high  
93 costs, long timelines, and extraordinarily high attrition rates in clinical trials, attributed  
94 to limited efficacy, safety concerns, off-target effects, or sometimes purely economic  
95 reasons (DiMasi *et al.* 2016; Minikel *et al.* 2024) Collectively, against this backdrop of  
96 failure, the need for transformative solutions for drug discovery becomes clear.  
97 Especially when we consider our incomplete understanding of target mechanism and  
98 the vast chemical space of compounds that can interact with that target.

## 99 100 **The role of AI in drug discovery**

101 Ideally, we would develop a comprehensive mathematical framework to systematically  
102 navigate the vast search spaces and intricate interactions inherent to drug discovery.  
103 However, realising such a framework has proven to be an immensely challenging  
104 endeavour with limited success so far. In contrast, methods using artificial intelligence  
105 (AI) are particularly well-suited for modelling the complexities and nuances of drug  
106 discovery. When employing AI, we essentially shift our approach: rather than relying  
107 on explicit mathematical descriptions of the underlying biology, chemistry, and  
108 physics, we leverage AI models to learn and infer patterns directly from data. While  
109 adopting data-driven machine learning techniques holds great promise for enhancing  
110 drug discovery pipelines, there are also certain trade-offs, such as a lack of  
111 transparency in the models and obscured understanding of causality.

112  
113 AI has potential to accelerate drug discovery by improving the identification of drug  
114 candidates and enhancing our understanding of their mechanisms. The increasing  
115 volume of diverse biological and chemical data, including genomics, proteomics,  
116 metabolomics, electronic health records, and biomedical literature, combined with  
117 high-throughput experiments, greatly enhances AI's ability to extract and interpret  
118 insights. Notably, recent studies have highlighted the importance of including genetic  
119 and genomic data in drug target discovery pipelines (Razuvayevskaya *et al.* 2023).  
120 One estimate quantifying the impact genetic evidence has on success of clinical trials,  
121 estimated the odds of advancing to a later stage of clinical trials to be 80% higher  
122 when genetic evidence for a target is present (Minikel *et al.* 2024). Furthermore, AI  
123 can be used to develop *in silico* methods to predict and simulate biological and  
124 chemical spaces. Examples of such approaches are cellular and genetic perturbation  
125 modelling (Bunne *et al.* 2023; Prasad *et al.* 2022), gene expression prediction (Avsec  
126 *et al.* 2021; Kelley *et al.* 2018; Linder *et al.* 2023), variant effect prediction (Brandes  
127 *et al.* 2022; Cheng *et al.* 2023; Frazer *et al.* 2021; Lin *et al.* 2023a), protein structure  
128 prediction (Baek *et al.* 2021; Jumper *et al.* 2021; Lin *et al.* 2023b), drug-target  
129 interaction prediction (Chen *et al.* 2016; Huang *et al.* 2021; Wen *et al.* 2017), and  
130 molecular docking simulations for drug design (Corso *et al.* 2023; Gentile *et al.* 2020).



131  
132 Figure 1: Venn diagram of guiding criteria for the maximum impact of AI in relation to  
133 drug discovery. We have made the connection to drug target discovery in the  
134 respective sets. The intersection of all sets is where the sweet spot for using AI lies.

135  
136 When it comes to determining the applicability of AI, we can refer to some guiding  
137 principles (Figure 1) that can help us to establish whether introducing AI to solve our  
138 problem is sensible. We argue that drug target discovery problems lie at the  
139 intersection of all these principles, making them amenable to be solved with AI.

140  
141 Firstly, the problem at hand must have sufficient scale. Building a successful AI model  
142 is reliant on having examples to learn from. While unsupervised approaches can be  
143 powerful, the potential of AI predominantly resides in the ability to uncover  
144 generalisable patterns within training data through a supervised or a self-supervised  
145 framework. A part of this scale is the quality of the data. The dataset should not just  
146 be large, but it should also be of  
147 high quality or be processed such that it is of high quality. High quality data implies  
148 that the model can learn meaningful signals from the patterns and relationships  
149 contained within the data. Some concrete examples of factors potentially decreasing  
150 data quality are noise, class imbalances, population bias, and missing data.

151  
152 Secondly, the complexity of the problem should be appropriate to fully leverage the  
153 power of AI models. At the lower bound of the complexity spectrum, the problem could  
154 be insufficiently complex, making it likely that an overparameterized AI model is  
155 developed that performs seemingly well, but does not generalise. This phenomenon

156 is referred to as overfitting in AI literature. Note that overfitting is not limited to this  
157 scenario and can also occur in poorly designed AI models where the problem itself is  
158 not necessarily insufficiently complex. At the other end of the complexity spectrum,  
159 a problem could be intractable. Take the entire chemical space of  $\sim 10^{60}$  compounds  
160 for example (Reymond 2015), this immense search space is simply too large for any  
161 computational method to fully explore. However, we can make this task more  
162 manageable by focusing on a smaller, more relevant subset of compounds. One  
163 effective approach to achieve this is by using generative AI models. These models are  
164 trained by adding random variations to existing, known data and then attempting to  
165 reconstruct the original input from this altered data. Through this process, the model  
166 learns the patterns and distributions inherent in the data which can be used to  
167 construct outputs based on these patterns.

168

169 In the context of drug discovery, this technique can be applied to known chemical  
170 structures. This is the basis of Generative Molecular Design (GMD), where AI models  
171 are used to generate potentially viable chemical compounds by learning from existing  
172 chemical structures (Thomas *et al.* 2023). This approach helps streamline the search  
173 for new drug candidates by focusing on the most promising areas of the vast chemical  
174 space, in this case up to  $\sim 10^{11}$  compounds (Ruddigkeit *et al.* 2012), constraining the  
175 search space and thus making the problem computationally tractable. For AI methods  
176 to thrive, a balance must be struck as it pertains to the complexity of the problem.  
177 We argue that drug discovery, including drug target discovery, satisfies the complexity  
178 criterion. Target discovery is often constrained to parameterisations of the genome,  
179 or the druggable genome. These are about 20,000 and 4,000 genes in size  
180 respectively, which is a tractable search space. As for the chemistry of compounds  
181 binding to the target, we can narrow down the search space to effectively design novel  
182 compounds.

183

184 Lastly, the input features for the problem should be non-linearly related to the target  
185 variable. Most biological phenomena are highly non-linear, so it is rare to encounter a  
186 biological problem where input and output are linearly related. This also becomes  
187 apparent from examining the AI models that underpin some seminal breakthroughs  
188 in the context of biology, such as CellOT for gene perturbation prediction (Bunne *et al.*  
189 *et al.* 2023), ESMFold and AlphaFold for protein structure prediction (Jumper *et al.* 2021;  
190 Lin *et al.* 2023b), and EVE and AlphaMissense for missense variant pathogenicity

191 prediction (Cheng *et al.* 2023; Frazer *et al.* 2021). To model the non-linearity inherent  
192 to these problems, non-linear activation functions are one of the key elements allowing  
193 AI models to effectively capture the highly complex relationships within the underlying  
194 distributions they attempt to model. Since many biological phenomena exhibit strong  
195 non-linearity, it makes sense to express and solve these problems in the language of  
196 AI.



## 197 **AI Methods and Data Modalities in Drug Target Discovery**

198 One leading reason for the convergence between AI and drug discovery is the diverse  
199 range of data types that are being used in drug discovery. The data can be presented  
200 in various forms, such as tabular, text, sequences, graphs, and images, each offering  
201 a distinct perspective into the biology underlying disease and potential cures. In Table  
202 1, we summarise the different modalities, their use-cases, and some open-access data  
203 sources. In the following paragraphs, we briefly discuss each data modality, and how  
204 it generally is used in drug target discovery.

205  
206 One of the most common methods for presenting data related to drug target discovery  
207 is through structured tables. Typically, these tabular data structures will contain  
208 information describing genes or variants, e.g., allele frequency, mutation type, and  
209 conservation scores across species. There are different resources and consortia that  
210 aggregate and characterise genomic data in tabular form, such as UK Biobank (Sudlow  
211 *et al.* 2015), Genes & Health (Finer *et al.* 2020), and Open Targets (Ochoa *et al.*  
212 2021). Traditional machine learning (ML) methods, e.g. XGBoost (Chen and Guestrin  
213 2016), Linear Regression, Logistic Regression (Pedregosa *et al.* 2011), as well as deep  
214 neural networks (LeCun *et al.* 2015), have been developed and tailored to tabular  
215 datasets. Therefore, these models have a track record of delivering outstanding  
216 performance when working with tabular data.

217  
218 Textual data, comprising scientific literature, research articles, patents, clinical trial  
219 reports, medical textbooks, chemical databases, and electronic health records,  
220 represents a valuable resource for drug discovery. The unstructured information in  
221 textual documents can provide us with critical insights related to potential drug  
222 targets, novel or repurposed drug candidates, and adverse events amongst others.  
223 Textual data is typically best analysed using Natural Language Processing (NLP)  
224 methods. Recently, Large Language Models (LLMs) have surfaced as the state-of-the-  
225 art model type to analyse textual data. LLMs are deep neural networks that combine  
226 many different layer types, such as embedding layers, attention layers and linear  
227 layers that coalesce to learn semantic information from textual input. Typically, LLMs  
228 are pre-trained using self-supervised approaches where a large corpus of text gets  
229 tokenised, i.e., it gets mapped to numerical vectors representing the words. This  
230 corpus is masked at random, and consequently tasked with predicting the next tokens  
231 (Devlin *et al.* 2019; Radford *et al.* 2018). For task-specific objectives, the pre-trained



232 model can be trained further on data related to the task of interest, e.g. information  
233 retrieval or translation. (Microsoft Research AI4Science and Microsoft Azure Quantum  
234 2023; Singhal *et al.* 2023a, 2023b)

235  
236 Data that can be represented sequentially are fundamental to biology. Such sequences  
237 often correspond to biological or chemical structure. Some of these data are genomic  
238 data, transcriptomic data, protein sequences, and drug compound libraries in the form  
239 of SMILES or SELFIES strings. Previously, we introduced language models within the  
240 context of natural language. Yet, their versatility transcends the domain of language.  
241 Language models also prove adept at understanding biological languages, e.g.,  
242 decoding semantic meaning from DNA via nucleotide sequences, and unravelling  
243 structural or functional information for proteins through the interpretation of amino  
244 acid sequences. To model and use these sequences, languages models can be trained  
245 to predict masked nucleotides or amino acids and consequently generalise to unseen  
246 sequences (Benegas *et al.* 2023; Dee 2022; Lin *et al.* 2023b). Another type of model  
247 showing promise on sequential and structural data are generative models. Generative  
248 models are self-supervised machine learning models that are trained to model the  
249 statistical distribution of input data, typically by reconstructing the original distribution  
250 after random noise has been added as input during the training process (Goodfellow  
251 *et al.* 2014). A couple of ways in which these models can be applied is to model DNA  
252 regulatory sequences (Zrimec *et al.* 2022), and they can be utilised to generate novel  
253 protein structures that meet some specified criteria. (Ingraham *et al.* 2023; Watson  
254 *et al.* 2023). Attention-based neural networks have shown to be well versed in  
255 analysing sequences to correct consensus sequence errors (Baid *et al.* 2023),  
256 comprehend protein structures (Baek *et al.* 2021; Lin *et al.* 2023b), and discover  
257 potential drug targets (Chen *et al.* 2023). The attention mechanism allows the model  
258 to learn relations between different parts of the input sequence, even if these parts  
259 are located far away from each other in their representation space (Vaswani *et al.*  
260 2017). The most notable example of an attention-based neural network working with  
261 sequence-based data is AlphaFold. AlphaFold predicts protein structure in 3D from an  
262 amino acid sequence input (Jumper *et al.* 2021).

263  
264 Network data (e.g., gene and protein interaction networks) can provide a  
265 comprehensive view of molecular relationships, by representing them efficiently as  
266 graphs with nodes and edges. Furthermore, representing data as a graph allows us to

267 build Graph Neural Networks (GNNs) (Veličković 2023). GNNs are optimised to learn  
268 and propagate information across nodes, allowing for efficient learning from these  
269 data structures. In the context of drug target discovery, there are various successful  
270 examples of graphs being used, such as in network expansion for pleiotropy mapping  
271 (Barrio-Hernandez *et al.* 2023), CausalBench (Chevalley *et al.* 2022), and many others  
272 (Muzio *et al.* 2021). A recent trend in drug target discovery has been the usage of  
273 Knowledge Graphs (KGs). These typically are heterogeneous graphs that store  
274 different data about compounds or genes in nodes, and relationships between nodes  
275 in the edges (Chandak *et al.* 2023).

276  
277 Medical imaging, including x-rays, CT scans, MRI, and histopathology slides, function  
278 as important assets for disease diagnosis and tracking treatment responses.  
279 Generative models, Convolutional Neural Networks (CNNs), Visual Transformers (ViTs)  
280 and deep learning architectures are frequently used for the analysis of visual data  
281 (Dosovitskiy *et al.* 2021; Liu *et al.* 2017; Tu *et al.* 2023). When it comes to molecular  
282 imaging, images are captured in various resolutions all the way down from the tissue  
283 to the cellular level. These images offer profound insights into the molecular intricacies  
284 of diseases and drug interactions. Finally, drug screening assays generate a treasure  
285 trove of image data, showcasing cells or organisms under perturbation of various  
286 compounds in pursuit of potential drugs. AI models help with their ability to  
287 comprehensively analyse the resulting images. Next to interpreting the images, using  
288 image data also often involves image correction and automatic feature extraction,  
289 both tasks in which AI methods excel (Dee *et al.* 2023; Krentzel *et al.* 2023).

290  
291 While it is true that certain data modalities conventionally have been associated with  
292 certain types of AI architectures, a lot of the state-of-the-art models do not exclusively  
293 use a single data modality or a single architecture. Often, data and model types are  
294 combined. This combination can occur in various ways and often different model types  
295 are involved with the processing of various types of data before it gets combined,  
296 which often happens in so-called embeddings (Alwazzan *et al.* 2023; Khader *et al.*  
297 2023; Ngiam *et al.* 2011; Venugopalan *et al.* 2021). Embeddings are representations  
298 of the raw input data in a latent space that can be used for downstream computations.  
299 Furthermore, most modern-day AI architectures consist of various blocks, which are  
300 organisational units in a neural network that are composed of different layers, or even

301 whole models that feed into each other and interact with each other. Models like this  
302 are often referred to as multimodal machine learning models.

	Data modality	Biological representation	Main AI architectures	Example data sources <sup>1</sup>
303				
304				
305	Tabular	Multiomics, Electronic Health Records	Traditional machine learning <sup>2</sup> , Multilayer perceptron	UK Biobank, Genes & Health, OpenTargets, TCGA, GEO
306				
307				
308				
309	Text	Gene ontology, Scientific literature, Clinical trials	Large language model	GO, PubMed, ClinicalTrials.gov
310				
311				
312				
313	Sequence/ Structure	DNA, RNA, Protein, Small molecules	Attention-based neural network, Generative model, Language model	Ensembl, UniProt, UCSC Genome Browser, ChEMBL, GenBank, PDB, GENCODE
314				
315				
316				
317				
318				
319				
320	Graph	PPI, Gene interaction network, Protein structures, Small molecule structures, Pathway annotations	Graph neural network	STRING, STITCH, BioGRID, PDB, TRRUST, RegNetwork, IntAct, PubChem, ChEMBL, Reactome, KEGG
321				
322				
323				
324				
325				
326				
327				
	Image	Histopathology, Radiology, Spatial transcriptomics	Convolutional neural network, Visual transformers	TCIA, GDC, MICA-MIC

328 Table 1: Categorisation of various data modalities commonly used in the field of  
 329 biomedical research and drug target discovery, along with biology the data represent,  
 330 the primary AI architecture employed on them, and key data sources. Note that the  
 331 AI architectures are not exclusive to these data modalities and in practise. Moreover,  
 332 often multiple are combined or sometimes even integrated into each other in an end-  
 333 to-end fashion.

334 <sup>1</sup> Citations to databases can be found in Supplementary Materials S2.

335 <sup>2</sup> In this case, we mean traditional machine learning to encompass linear and logistic regression, support vector  
336 machines and tree boosting models.

### 337 **Exploring AI-Based Strategies for Drug Target Identification**

338 The first example we will explore is DrugnomeAI, an ensemble architecture for the  
339 prediction of drug targets (Polikar 2006; Raies *et al.* 2022; Vitsios and Petrovski  
340 2020). DrugnomeAI excels in predicting the druggability of candidate drug targets by  
341 leveraging 324 gene-level features for every protein-coding gene within the human  
342 exome. Raies *et al.* conducted a feature importance study with Boruta, which is a  
343 feature selection technique that helps identify the most relevant variables in a dataset  
344 by comparing their importance to that of randomised, noise-added variables (Kursa  
345 *et al.* 2010). This analysis showed that the most informative features for druggability  
346 prediction were protein-protein interaction features. This is in line with existing  
347 research showing that partners of druggable genes are also likely to be druggable  
348 (Finan *et al.* 2017). Raies and colleagues frame their model's objective as a positive-  
349 unlabelled learning (PUL) problem. Here, the positive dataset comprises targets for  
350 which they have identified evidence of druggability, while the unlabelled set  
351 encompasses the remaining targets. The ultimate task is to rank these remaining  
352 targets based on their predicted druggability. Within their PUL framework, Raies *et al.*  
353 use eight separate classifiers that are stochastically trained through a 10-fold cross-  
354 validation process. Subsequently, the predictions from these classifiers are combined  
355 to produce the final ranking of the unlabelled drug targets. Notably, Raies *et al.*  
356 observed that the top-ranked genes in their prioritisation exhibit significant  
357 enrichment in the clinical literature, arguing that their model has effectively recognised  
358 druggability patterns within the feature set.

359  
360 It is also possible to combine multiple data modalities in a more direct way than  
361 ensemble modelling, namely via multitask learning (Caruana 1998). A multitask  
362 learning problem in drug target discovery is typically framed as one where you are  
363 trying to predict target qualities as well as properties of the target-binding drug (Lin  
364 *et al.* 2022; Sadawi *et al.* 2019). Multitask learning allows the model to co-learn a set  
365 of tasks together to optimise overall performance. This approach leverages shared  
366 information between tasks, combatting overfitting and improving generalisation.  
367 Multitask neural networks can integrate data from various sources, making them  
368 valuable for a wide range of tasks such as predicting drug targets, but also drug  
369 toxicity and sensitivity (Ammad-Ud-Din *et al.* 2017; Costello *et al.* 2014).  
370 Furthermore, they offer a means to bridge the gap between biology and chemistry in

371 drug discovery by incorporating structural data like SMILES representations, next to  
372 information characterising the biological target, enabling simultaneous prediction of  
373 side effects, ligand docking, likely targets, and key compound properties (Mikolov *et*  
374 *al.* 2013b, 2013a).

375 In some areas of study where data is sparsely available, such as for rare diseases or  
376 diseases in clinically unavailable tissues, AI methods can meaningfully identify  
377 candidate drug targets through transfer learning. Transfer learning is a concept in AI  
378 where we train on abundant data that is tangentially related to some problem with  
379 limited data, and consequently fine-tune the resulting model towards the limited data  
380 case (Pan and Yang 2010). One example of a model utilising transfer learning is  
381 Geneformer (Theodoris *et al.* 2023). Geneformer uses self-attention to pick out  
382 important genes using transcriptomic data, which can vary across different cell types,  
383 developmental stages, or disease conditions. Geneformer was trained with a dataset  
384 called Genecorpus-30M, which was assembled from 29.9 million human single-cell  
385 transcriptomes. The transcriptome data is processed through six transformer encoder  
386 units, involving self-attention and feed-forward layers. Pre-training is done using a  
387 masked learning objective, where 15% of genes in each transcriptome are masked,  
388 and the model learns to predict the masked genes based on the context of the  
389 unmasked genes. Due to the size and broad scope of Geneformer's pre-training,  
390 together with the potential to fine-tune the model, we refer to this model as a  
391 foundation model (Bommasani *et al.* 2022). Using Geneformer, cardiomyocytes from  
392 three types of limitedly available heart tissue were studied: healthy (n=9),  
393 hypertrophic cardiomyopathy (n=11), or dilated cardiomyopathy (n=9). Theodoris *et*  
394 *al.* performed *in silico* treatment analysis by either inhibiting or activating pathways  
395 and seeing if this would move the healthy cell states towards either hypertrophic or  
396 dilated cardiomyopathic states. If so, the pathway was inspected for potential  
397 therapeutic targets based on druggability and disease relevance. A target that was  
398 highlighted through this analysis was *ADCY5*, which is a known druggable target  
399 (Wagner *et al.* 2016a) as well as involved in longevity and protection of  
400 cardiomyocytes in mouse models (Ho *et al.* 2010). Another target that *in silico*  
401 treatment analysis pointed to in this context was *SRPK3*, which is a downstream  
402 effector of *MEF2* (Nakagawa *et al.* 2005). *MEF2* is known to play a role in myocardial  
403 cell hypertrophy (Akazawa and Komuro 2003). While single-cell foundation models  
404 have demonstrated impressive results in certain situations and seem conceptually  
405 attractive for downstream applications, it's important to exercise caution. These pre-

406 trained models may not perform well in all contexts, particularly for zero-shot  
407 prediction in other biological contexts (Kedzierska *et al.* 2023). Therefore, employing  
408 biological foundation models for zero-shot prediction in contexts divergent from their  
409 original training objective should be approached carefully.

410  
411 GNNs are also being employed in drug target discovery. One such approach is EMOGI  
412 (Schulte-Sasse *et al.* 2021), a graph convolutional network (GCN) that predicts cancer  
413 drug targets. EMOGI stands out by integrating a wide range of interaction and  
414 multiomics data to predict cancer genes. This way of combining different data sources  
415 addresses the evolving understanding of cancer as a complex interplay of genetic and  
416 non-genetic factors (Bell and Gilan 2020; Hanahan and Weinberg 2011). Unlike  
417 previous approaches that primarily rely on somatic mutations and occasionally  
418 integrate PPI networks (Cowen *et al.* 2017; Leiserson *et al.* 2015; Reyna *et al.* 2018),  
419 EMOGI employs GCNs to predict cancer genes by amalgamating multiple data  
420 modalities, including mutations, copy number variations, DNA methylation, gene  
421 expression, and PPI networks. The graph is constructed to have its topology represent  
422 a PPI network. This means that the nodes represent genes, and the edges represent  
423 whether two genes interact. R. Schulte-Sassen *et al.* also did interpretability analysis  
424 of their GCN model. They use the Layer-wise Relevance Propagation (LRP) propagation  
425 rule (Bach *et al.* 2015), which allows for dissecting what is happening in the GCNs and  
426 provides us with insights into why specific genes are classified as cancer-related.  
427 Through biclustering and LRP analysis, distinct modules of newly predicted cancer  
428 genes (NPCGs) are revealed—some predominantly influenced by network interactions,  
429 others primarily driven by omics features. These NPCGs, while not always necessarily  
430 displaying recurrent alterations themselves, interact with known cancer genes,  
431 positioning them as significant players in tumorigenesis. Notably, these predictions  
432 align with essential genes identified through loss-of-function screens, reinforcing the  
433 credibility of EMOGI's insights.

434  
435 Beyond academic research and applications, as of Q3 2023, there are a plethora of  
436 AI-derived therapeutics in clinical trial pipelines. Most of these come forth out of  
437 industrial research laboratories. A lot of the information that is publicly available on  
438 how AI is influencing drug target discovery comes from what we here refer to as AI-  
439 first drug discovery companies. These are companies that highlight explicitly the fact  
440 that they are using AI in their drug target discovery and drug design efforts. While we



441 can only associate drugs being AI-derived from such companies, we should note that  
442 big pharmaceutical companies are also heavily investing into introducing AI into their  
443 pipelines. However, it is much harder to attribute the involvement of AI in the  
444 development of new pharmaceuticals in this case. So, while looking at the status of  
445 AI-first companies might be a good probe into the penetrance of AI into the  
446 pharmaceutical industry, it does not provide us with a comprehensive view of the role  
447 AI is currently playing in industry.

448

449 In Figure 2, we have visualised the status of targets and associated compounds  
450 currently in clinical and preclinical trials. The data was put together by searching and  
451 collecting a list of publicly and privately held companies that explicitly mention the  
452 usage of AI on their website. We have added a table containing the data we collected  
453 in Supplementary Table S1. Note that this is not an exhaustive list and we only  
454 included target-compound pairs for which we could find sufficient data in the pipelines  
455 reported by the companies. For discontinued compounds, press-releases and historical  
456 website snapshots have been consulted to confirm the development status of  
457 compounds. The discontinued compounds collected in our data is an underestimation  
458 of the true number of discontinued compounds. Often, data and status on discontinued  
459 compounds is not easily accessible in public records. Hence, the only discontinued  
460 compounds added in this list, are ones that (i) have had accessible press coverage,  
461 (ii) have been withdrawn from a clinical trial investigation as indicated by  
462 ClinicalTrials.gov, or (iii) have been mentioned in an accessible snapshot of a  
463 company's pipeline webpage, consulted via wayback.archive.org, and removed  
464 without any mention of success. We only consider compounds in which the company  
465 was leading the effort for approval. We use FDA approval status to determine whether  
466 a compound has been officially approved. We excluded AI-first companies that have  
467 not yet had at least one compound enter clinical trials.

468

469

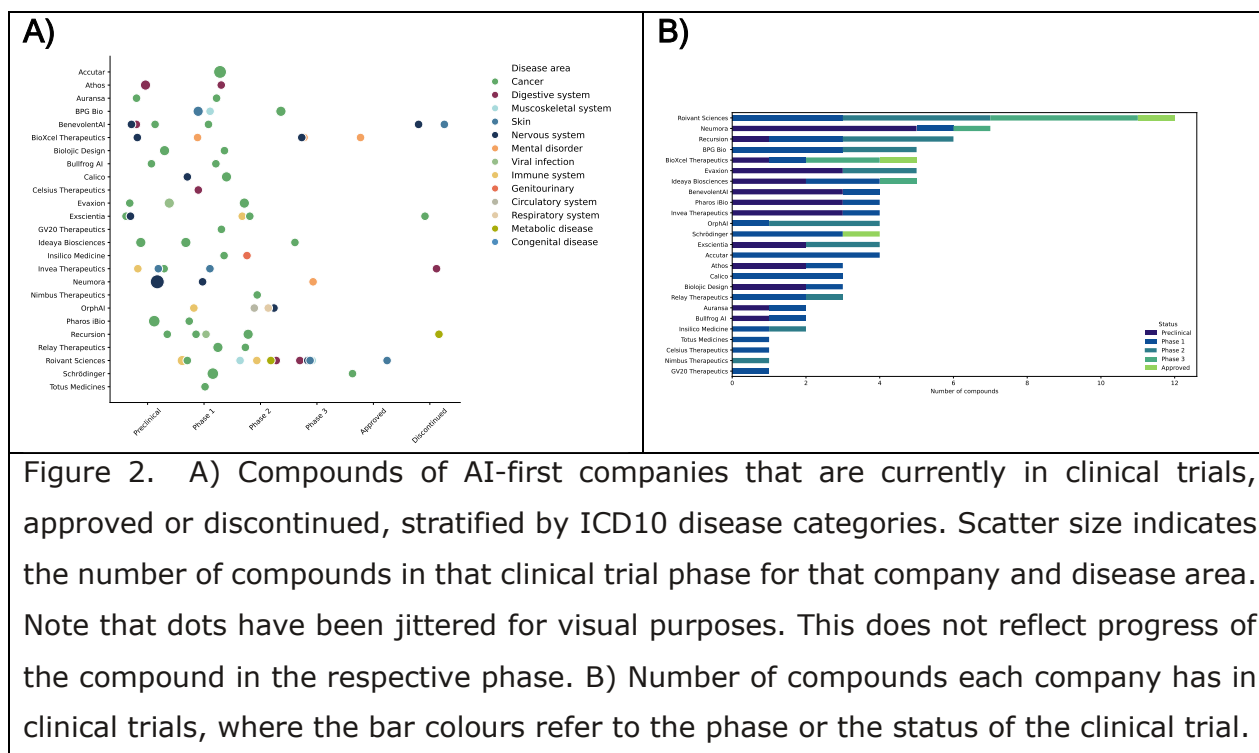


Figure 2. A) Compounds of AI-first companies that are currently in clinical trials, approved or discontinued, stratified by ICD10 disease categories. Scatter size indicates the number of compounds in that clinical trial phase for that company and disease area. Note that dots have been jittered for visual purposes. This does not reflect progress of the compound in the respective phase. B) Number of compounds each company has in clinical trials, where the bar colours refer to the phase or the status of the clinical trial.

470

## 471 Discussion and Future Prospects

472 AI is penetrating all levels of drug discovery, including target discovery and validation.  
 473 AI methods rely on the existence of large, high quality data sets. Currently, these data  
 474 exist but are certainly incomplete and potentially confounding in nature. We must take  
 475 note of the limitations of existing data and look at ways to improve data in a targeted  
 476 manner. Most publicly available big data sets often rely on aggregated information  
 477 descendent from skewed representations of the population. Different populations  
 478 display widely varying genomic characteristics and responses to drugs, and  
 479 consequently, less represented populations suffer from diminished treatment  
 480 outcomes (Gross *et al.* 2022; Popejoy and Fullerton 2016; Ramamoorthy *et al.* 2015).  
 481 Therefore, the databases used to identify drug targets often lack sufficient  
 482 representation of population diversity, resulting in disparate health outcomes for  
 483 diseases that are effectively treated in well-represented groups but remain challenging  
 484 to address in the underrepresented populations. (Hindorff *et al.* 2018; Landry *et al.*  
 485 2018).

486

487 At the molecular level, we encounter a different set of biases in the data we use to  
 488 train our models. For example, some protein classes are significantly overrepresented

489 compared to others based on FDA approval data, which may be attributed to shared  
490 structural or functional similarities for proteins within a given class. If we train a new  
491 generation of models with these targets as labels, we are likely to perpetuate these  
492 biases in newly prioritised drug targets. Furthermore, we should also acknowledge  
493 that because of data availability limitations, bias and historical momentum around  
494 known drug targets and classes of targets, there is a significant portion of the genome  
495 of which we know too little to assess their validity as drug targets (Finan *et al.* 2017;  
496 Oprea *et al.* 2018; Wood *et al.* 2019). Assuming there are also potential drug targets  
497 hidden within what has been colloquially termed the "unknome" (Rocha *et al.* 2023),  
498 this would increase the search space of potential drug targets further beyond what  
499 the current paradigm of what drug target druggability models consider. Another  
500 challenge is that the concept of a druggable target is not static. This is particularly  
501 pronounced for cancer, where target-associated pathways are prone to quickly  
502 becoming resistant to treatment through various mechanisms (Shabani and Hojjat-  
503 Farsangi 2016). This means that the "one disease, one target" paradigm might not be  
504 the best approach to curing diseases, even in cases where a single target is indeed  
505 initially therapeutically receptive to treat the disease.

506  
507 While AI-powered drug target discovery has its fair share of obstacles to overcome, it  
508 is still a field that is in its infancy. Moreover, next to these obstacles lie many  
509 opportunities for promising discoveries. This is not only limited to drug target  
510 discovery, but drug discovery in its broadest sense. For the successful application of  
511 AI, specifically deep learning-based architectures, the three guiding principles must  
512 be satisfied: scale, complexity, and non-linearity. We argue that drug target discovery  
513 satisfies all three of these principles. Given this reality, AI-based methods stand to  
514 improve the speed with which we can discover and validate novel drug targets. Recent  
515 breakthroughs in AI have led to improvements by providing an increased ability to  
516 incorporate sequence and structure-based target evidence. As models like AlphaFold  
517 are improved and extended to also reflect the dynamic nature of proteins, and we  
518 incorporate small molecules and macromolecular structures into these models, our  
519 ability to do *in silico* drug discovery will dramatically improve. In addition to predicting  
520 protein structures, AI methods stand to significantly improve a multitude of other  
521 biological challenges. These include, but are not limited to, predicting gene  
522 perturbations, assessing the effects of genetic variants, *de novo* generation of  
523 proteins, and molecular docking simulations. In the long run, transitioning a significant

524 portion of the drug discovery pipeline to an *in silico* environment holds substantial  
525 advantages for all parties involved with drug discovery. For patients, this shift would  
526 enhance the efficiency of developing new and safe medications, resulting in faster  
527 delivery of improved therapeutics. For pharmaceutical companies, this transition  
528 would lead to significant cost and time savings, which are estimated between 25%  
529 and 50% up to the preclinical stage (Loynachan *et al.* 2023). For us to get to this  
530 point, experimental validations of *in silico* methods remain essential both to validate  
531 computational predictions and to provide labels for the models to train with.

532

533 AI-driven drug target discovery presents a promising avenue for identifying novel,  
534 safe and efficacious targets. By leveraging the abundance of multiomics data, and the  
535 power of modern AI architectures, applicable to a variety of data modalities – ranging  
536 from images to sequences and protein structures – we find ourselves at the precipice  
537 of having data and method converge at meaningful impact on drug target discovery,  
538 and drug discovery at large.

539

540 **References**

- 541 **Abdel-Magid AF** (2015) Allosteric Modulators: An Emerging Concept in Drug  
542 Discovery. *ACS Medicinal Chemistry Letters* **6**(2), 104-107.  
543 <https://doi.org/10.1021/ml5005365>.
- 544 **Akazawa H and Komuro I** (2003) Roles of Cardiac Transcription Factors in Cardiac  
545 Hypertrophy. *Circulation Research* **92**(10), 1079-1088.  
546 <https://doi.org/10.1161/01.RES.0000072977.86706.23>.
- 547 **Alwazzan O, Khan A, Patras I and Slabaugh G** (2023) MOAB: Multi-Modal Outer  
548 Arithmetic Block for Fusion of Histopathological Images and Genetic Data for  
549 Brain Tumor Grading. In *2023 IEEE 20th International Symposium on*  
550 *Biomedical Imaging (ISBI)*. 1-5.  
551 <https://doi.org/10.1109/ISBI53787.2023.10230698>.
- 552 **Ammad-Ud-Din M, Khan SA, Wennerberg K and Aittokallio T** (2017) Systematic  
553 identification of feature combinations for predicting drug response with  
554 Bayesian multi-view multi-task linear regression. *Bioinformatics (Oxford,*  
555 *England)* **33**(14), i359-i368. <https://doi.org/10.1093/bioinformatics/btx266>.
- 556 **Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski**  
557 **K, Dwight SS, Eppig JT, Harris MA, Hill DP, Issel-Tarver L, Kasarskis A, Lewis**  
558 **S, Matese JC, Richardson JE, Ringwald M, Rubin GM and Sherlock G** (2000)  
559 Gene Ontology: tool for the unification of biology. *Nature Genetics* **25**(1), 25-  
560 29. <https://doi.org/10.1038/75556>.
- 561 **Avsec Ž, Agarwal V, Visentin D, Ledsam JR, Grabska-Barwinska A, Taylor KR,**  
562 **Assael Y, Jumper J, Kohli P and Kelley DR** (2021) Effective gene expression  
563 prediction from sequence by integrating long-range interactions. *Nature*  
564 *Methods* **18**(10), 1196-1203. <https://doi.org/10.1038/s41592-021-01252-x>.

- 565 **Bach S, Binder A, Montavon G, Klauschen F, Müller K-R and Samek W (2015)** On  
566 Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise  
567 Relevance Propagation. *PLOS ONE* **10**(7), e0130140.  
568 <https://doi.org/10.1371/journal.pone.0130140>.
- 569 **Baek M, DiMaio F, Anishchenko I, Dauparas J, Ovchinnikov S, Lee GR, Wang J,**  
570 **Cong Q, Kinch LN, Schaeffer RD, Millán C, Park H, Adams C, Glassman CR,**  
571 **DeGiovanni A, Pereira JH, Rodrigues AV, van Dijk AA, Ebrecht AC,**  
572 **Opperman DJ, Sagmeister T, Buhlheller C, Pavkov-Keller T, Rathinaswamy**  
573 **MK, Dalwadi U, Yip CK, Burke JE, Garcia KC, Grishin NV, Adams PD, Read**  
574 **RJ and Baker D (2021)** Accurate prediction of protein structures and  
575 interactions using a three-track neural network. *Science* **373**(6557), 871-876.  
576 <https://doi.org/10.1126/science.abj8754>.
- 577 **Baid G, Cook DE, Shafin K, Yun T, Llinares-López F, Berthet Q, Belyaeva A, Töpfer**  
578 **A, Wenger AM, Rowell WJ, Yang H, Kolesnikov A, Ammar W, Vert J-P,**  
579 **Vaswani A, McLean CY, Nattestad M, Chang P-C and Carroll A (2023)**  
580 DeepConsensus improves the accuracy of sequences with a gap-aware  
581 sequence transformer. *Nature Biotechnology* **41**(2), 232-238.  
582 <https://doi.org/10.1038/s41587-022-01435-7>.
- 583 **Bakheet TM and Doig AJ (2009)** Properties and identification of human protein drug  
584 targets. *Bioinformatics* **25**(4), 451-457.  
585 <https://doi.org/10.1093/bioinformatics/btp002>.
- 586 **Barrio-Hernandez I, Schwartzentruber J, Shrivastava A, del-Toro N, Gonzalez A,**  
587 **Zhang Q, Mountjoy E, Suveges D, Ochoa D, Ghousaini M, Bradley G,**  
588 **Hermjakob H, Orchard S, Dunham I, Anderson CA, Porras P and Beltrao P**  
589 **(2023)** Network expansion of genetic associations defines a pleiotropy map of

590 human cell biology. *Nature Genetics* 1-10. [https://doi.org/10.1038/s41588-](https://doi.org/10.1038/s41588-023-01327-9)  
591 023-01327-9.

592 **Bell CC and Gilan O** (2020) Principles and mechanisms of non-genetic resistance in  
593 cancer. *British Journal of Cancer* **122**(4), 465-472.  
594 <https://doi.org/10.1038/s41416-019-0648-6>.

595 **Benegas G, Batra SS and Song YS** (2023) DNA language models are powerful  
596 predictors of genome-wide variant effects. *Proceedings of the National*  
597 *Academy of Sciences* **120**(44), e2311219120.  
598 <https://doi.org/10.1073/pnas.2311219120>.

599 **Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN**  
600 **and Bourne PE** (2000) The Protein Data Bank. *Nucleic Acids Research* **28**(1),  
601 235-242. <https://doi.org/10.1093/nar/28.1.235>.

602 **Bommasani R, Hudson DA, Adeli E, Altman R, Arora S, von Arx S, Bernstein MS,**  
603 **Bohg J, Bosselut A, Brunskill E, Brynjolfsson E, Buch S, Card D, Castellon R,**  
604 **Chatterji N, Chen A, Creel K, Davis JQ, Demszky D, Donahue C, Doumbouya**  
605 **M, Durmus E, Ermon S, Etchemendy J, Ethayarajh K, Fei-Fei L, Finn C, Gale**  
606 **T, Gillespie L, Goel K, Goodman N, Grossman S, Guha N, Hashimoto T,**  
607 **Henderson P, Hewitt J, Ho DE, Hong J, Hsu K, Huang J, Icard T, Jain S,**  
608 **Jurafsky D, Kalluri P, Karamcheti S, Keeling G, Khani F, Khattab O, Koh PW,**  
609 **Krass M, Krishna R, Kuditipudi R, Kumar A, Ladhak F, Lee M, Lee T,**  
610 **Leskovec J, Levent I, Li XL, Li X, Ma T, Malik A, Manning CD, Mirchandani S,**  
611 **Mitchell E, Munyikwa Z, Nair S, Narayan A, Narayanan D, Newman B, Nie A,**  
612 **Niebles JC, Nilforoshan H, Nyarko J, Ogut G, Orr L, Papadimitriou I, Park JS,**  
613 **Piech C, Portelance E, Potts C, Raghunathan A, Reich R, Ren H, Rong F,**  
614 **Roohani Y, Ruiz C, Ryan J, Ré C, Sadigh D, Sagawa S, Santhanam K, Shih**



- 615 A, Srinivasan K, Tamkin A, Taori R, Thomas AW, Tramèr F, Wang RE, Wang  
616 W, Wu B, Wu J, Wu Y, Xie SM, Yasunaga M, You J, Zaharia M, Zhang M,  
617 Zhang T, Zhang X, Zhang Y, Zheng L, Zhou K and Liang P (2022, July 12) On  
618 the Opportunities and Risks of Foundation Models. arXiv.  
619 <https://doi.org/10.48550/arXiv.2108.07258>.
- 620 **Brandes N, Goldman G, Wang CH, Ye CJ and Ntranos V** (2022, August 26)  
621 Genome-wide prediction of disease variants with a deep protein language  
622 model. bioRxiv, 2022.08.25.505311.  
623 <https://doi.org/10.1101/2022.08.25.505311>.
- 624 **Bunne C, Stark SG, Gut G, del Castillo JS, Levesque M, Lehmann K-V, Pelkmans L,**  
625 **Krause A and Rätsch G** (2023) Learning single-cell perturbation responses  
626 using neural optimal transport. *Nature Methods* 1-10.  
627 <https://doi.org/10.1038/s41592-023-01969-x>.
- 628 **Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, Motyer A, Vukcevic**  
629 **D, Delaneau O, O'Connell J, Cortes A, Welsh S, Young A, Effingham M,**  
630 **McVean G, Leslie S, Allen N, Donnelly P and Marchini J** (2018) The UK  
631 Biobank resource with deep phenotyping and genomic data. *Nature*  
632 **562(7726)**, 203-209. <https://doi.org/10.1038/s41586-018-0579-z>.
- 633 **Caruana R** (1998) Multitask Learning. In Thrun S and Pratt L (eds), *Learning to*  
634 *Learn*. Boston, MA: Springer US, 95-133. [https://doi.org/10.1007/978-1-4615-](https://doi.org/10.1007/978-1-4615-5529-2_5)  
635 [5529-2\\_5](https://doi.org/10.1007/978-1-4615-5529-2_5).
- 636 **Chandak P, Huang K and Zitnik M** (2023) Building a knowledge graph to enable  
637 precision medicine. *Scientific Data* **10(1)**, 67. [https://doi.org/10.1038/s41597-](https://doi.org/10.1038/s41597-023-01960-3)  
638 [023-01960-3](https://doi.org/10.1038/s41597-023-01960-3).

- 639 **Chen J, Gu Z, Xu Y, Deng M, Lai L and Pei J** (2023) QuoteTarget: A sequence-  
640 based transformer protein language model to identify potentially druggable  
641 protein targets. *Protein Science* **32**(2), e4555.  
642 <https://doi.org/10.1002/pro.4555>.
- 643 **Chen T and Guestrin C** (2016) XGBoost: A Scalable Tree Boosting System. In  
644 *Proceedings of the 22nd ACM SIGKDD International Conference on*  
645 *Knowledge Discovery and Data Mining*. New York, NY, USA: Association for  
646 Computing Machinery, 785-794. <https://doi.org/10.1145/2939672.2939785>.
- 647 **Chen X, Yan CC, Zhang X, Zhang X, Dai F, Yin J and Zhang Y** (2016) Drug-target  
648 interaction prediction: databases, web servers and computational models.  
649 *Briefings in Bioinformatics* **17**(4), 696-712. <https://doi.org/10.1093/bib/bbv066>.
- 650 **Cheng J, Novati G, Pan J, Bycroft C, Žemgulytė A, Applebaum T, Pritzel A, Wong**  
651 **LH, Zielinski M, Sargeant T, Schneider RG, Senior AW, Jumper J, Hassabis**  
652 **D, Kohli P and Avsec Ž** (2023) Accurate proteome-wide missense variant  
653 effect prediction with AlphaMissense. *Science* **381**(6664), eadg7492.  
654 <https://doi.org/10.1126/science.adg7492>.
- 655 **Chevalley M, Roohani Y, Mehrjou A, Leskovec J and Schwab P** (2022, October 31)  
656 CausalBench: A Large-scale Benchmark for Network Inference from Single-  
657 cell Perturbation Data. arXiv. <http://arxiv.org/abs/2210.17283> (accessed 8  
658 March 2023)
- 659 **Clark K, Vendt B, Smith K, Freymann J, Kirby J, Koppel P, Moore S, Phillips S,**  
660 **Maffitt D, Pringle M, Tarbox L and Prior F** (2013) The Cancer Imaging Archive  
661 (TCIA): maintaining and operating a public information repository. *Journal of*  
662 *Digital Imaging* **26**(6), 1045-1057. <https://doi.org/10.1007/s10278-013-9622-7>.

- 663 **Clough E and Barrett T** (2016) The Gene Expression Omnibus database. *Methods in*  
664 *Molecular Biology (Clifton, N.J.)* **1418**, 93-110. [https://doi.org/10.1007/978-1-](https://doi.org/10.1007/978-1-4939-3578-9_5)  
665 [4939-3578-9\\_5](https://doi.org/10.1007/978-1-4939-3578-9_5).
- 666 **Corso G, Stärk H, Jing B, Barzilay R and Jaakkola T** (2023, February 11) DiffDock:  
667 Diffusion Steps, Twists, and Turns for Molecular Docking. arXiv.  
668 <http://arxiv.org/abs/2210.01776> (accessed 30 August 2023)
- 669 **Costello JC, Heiser LM, Georgii E, Gönen M, Menden MP, Wang NJ, Bansal M,**  
670 **Ammad-ud-din M, Hintsanen P, Khan SA, Mpindi J-P, Kallioniemi O, Honkela**  
671 **A, Aittokallio T, Wennerberg K, Collins JJ, Gallahan D, Singer D, Saez-**  
672 **Rodriguez J, Kaski S, Gray JW and Stolovitzky G** (2014) A community effort to  
673 assess and improve drug sensitivity prediction algorithms. *Nature*  
674 *Biotechnology* **32**(12), 1202-1212. <https://doi.org/10.1038/nbt.2877>.
- 675 **Cowen L, Ideker T, Raphael BJ and Sharan R** (2017) Network propagation: a  
676 universal amplifier of genetic associations. *Nature Reviews. Genetics* **18**(9),  
677 551-562. <https://doi.org/10.1038/nrg.2017.38>.
- 678 **Dee W** (2022) LMPred: predicting antimicrobial peptides using pre-trained language  
679 models and deep learning. *Bioinformatics Advances* **2**(1).  
680 <https://doi.org/10.1093/bioadv/vbac021>.
- 681 **Dee W, Sequeira I, Lobley A and Slabaugh G** (2023, December 13) Cell-Vision  
682 Fusion: A Swin Transformer-based Approach to Predicting Kinase Inhibitor  
683 Mechanism of Action from Cell Painting Data. bioRxiv, 2023.12.13.571534.  
684 <https://doi.org/10.1101/2023.12.13.571534>.
- 685 **Devlin J, Chang M-W, Lee K and Toutanova K** (2019, May 24) BERT: Pre-training of  
686 Deep Bidirectional Transformers for Language Understanding. arXiv.  
687 <https://doi.org/10.48550/arXiv.1810.04805>.

- 688 **DiMasi JA, Grabowski HG and Hansen RW (2016)** Innovation in the pharmaceutical  
689 industry: New estimates of R&D costs. *Journal of Health Economics* **47**, 20-  
690 33. <https://doi.org/10.1016/j.jhealeco.2016.01.012>.
- 691 **Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T,**  
692 **Dehghani M, Minderer M, Heigold G, Gelly S, Uszkoreit J and Hounsby N**  
693 (2021, June 3) An Image is Worth 16x16 Words: Transformers for Image  
694 Recognition at Scale. arXiv. <https://doi.org/10.48550/arXiv.2010.11929>.
- 695 **Drews J (2000)** Drug Discovery: A Historical Perspective. **287**(5460), 1960-1964.  
696 <https://doi.org/10.1126/science.287.5460.1960>.
- 697 **Finan C, Gaulton A, Kruger FA, Lumbers RT, Shah T, Engmann J, Galver L, Kelley**  
698 **R, Karlsson A, Santos R, Overington JP, Hingorani AD and Casas JP (2017)**  
699 The druggable genome and support for target identification and validation in  
700 drug development. *Science Translational Medicine* **9**(383), eaag1166.  
701 <https://doi.org/10.1126/scitranslmed.aag1166>.
- 702 **Finer S, Martin HC, Khan A, Hunt KA, MacLaughlin B, Ahmed Z, Ashcroft R, Durham**  
703 **C, MacArthur DG, McCarthy MI, Robson J, Trivedi B, Griffiths C, Wright J,**  
704 **Trembath RC and van Heel DA (2020)** Cohort Profile: East London Genes &  
705 Health (ELGH), a community-based population genomics and health study in  
706 British Bangladeshi and British Pakistani people. *International Journal of*  
707 *Epidemiology* **49**(1), 20-21i. <https://doi.org/10.1093/ije/dyz174>.
- 708 **Frankish A, Diekhans M, Jungreis I, Lagarde J, Loveland JE, Mudge JM, Sisu C,**  
709 **Wright JC, Armstrong J, Barnes I, Berry A, Bignell A, Boix C, Carbonell Sala**  
710 **S, Cunningham F, Di Domenico T, Donaldson S, Fiddes IT, García Girón C,**  
711 **Gonzalez JM, Grego T, Hardy M, Hourlier T, Howe KL, Hunt T, Izuogu OG,**  
712 **Johnson R, Martin FJ, Martínez L, Mohanan S, Muir P, Navarro FCP, Parker**

- 713 A, Pei B, Pozo F, Riera FC, Ruffier M, Schmitt BM, Stapleton E, Suner M-M,  
714 Sycheva I, Uszczyńska-Ratajczak B, Wolf MY, Xu J, Yang YT, Yates A,  
715 Zerbino D, Zhang Y, Choudhary JS, Gerstein M, Guigó R, Hubbard TJP, Kellis  
716 M, Paten B, Tress ML and Flicek P (2021) GENCODE 2021. *Nucleic Acids  
717 Research* 49(D1), D916-D923. <https://doi.org/10.1093/nar/gkaa1087>.
- 718 Frazer J, Notin P, Dias M, Gomez A, Min JK, Brock K, Gal Y and Marks DS (2021)  
719 Disease variant prediction with deep generative models of evolutionary data.  
720 *Nature* 599(7883), 91-95. <https://doi.org/10.1038/s41586-021-04043-8>.
- 721 Gentile F, Agrawal V, Hsing M, Ton A-T, Ban F, Norinder U, Gleave ME and  
722 Cherkasov A (2020) Deep Docking: A Deep Learning Platform for  
723 Augmentation of Structure Based Drug Discovery. *ACS Central Science* 6(6),  
724 939-949. <https://doi.org/10.1021/acscentsci.0c00229>.
- 725 Gillespie M, Jassal B, Stephan R, Milacic M, Rothfels K, Senff-Ribeiro A, Griss J,  
726 Sevilla C, Matthews L, Gong C, Deng C, Varusai T, Ragueneau E, Haider Y,  
727 May B, Shamovsky V, Weiser J, Brunson T, Sanati N, Beckman L, Shao X,  
728 Fabregat A, Sidiropoulos K, Murillo J, Viteri G, Cook J, Shorser S, Bader G,  
729 Demir E, Sander C, Haw R, Wu G, Stein L, Hermjakob H and D'Eustachio P  
730 (2022) The reactome pathway knowledgebase 2022. *Nucleic Acids Research*  
731 50(D1), D687-D692. <https://doi.org/10.1093/nar/gkab1028>.
- 732 Goodfellow IJ, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville  
733 A and Bengio Y (2014, June 10) Generative Adversarial Networks. arXiv.  
734 <https://doi.org/10.48550/arXiv.1406.2661>.
- 735 Gross AS, Harry AC, Clifton CS and Della Pasqua O (2022) Clinical trial diversity: An  
736 opportunity for improved insight into the determinants of variability in drug

- 737 response. *British Journal of Clinical Pharmacology* **88**(6), 2700-2717.  
738 <https://doi.org/10.1111/bcp.15242>.
- 739 **Grossman RL, Heath AP, Ferretti V, Varmus HE, Lowy DR, Kibbe WA and Staudt**  
740 **LM** (2016) Toward a Shared Vision for Cancer Genomic Data. *The New*  
741 *England Journal of Medicine* **375**(12), 1109-1112.  
742 <https://doi.org/10.1056/NEJMp1607591>.
- 743 **Han H, Cho J-W, Lee S, Yun A, Kim H, Bae D, Yang S, Kim CY, Lee M, Kim E, Lee**  
744 **S, Kang B, Jeong D, Kim Y, Jeon H-N, Jung H, Nam S, Chung M, Kim J-H**  
745 **and Lee I** (2018) TRRUST v2: an expanded reference database of human and  
746 mouse transcriptional regulatory interactions. *Nucleic Acids Research* **46**(D1),  
747 D380-D386. <https://doi.org/10.1093/nar/gkx1013>.
- 748 **Hanahan D and Weinberg RA** (2011) Hallmarks of Cancer: The Next Generation.  
749 *Cell* **144**(5), 646-674. <https://doi.org/10.1016/j.cell.2011.02.013>.
- 750 **Hann MM and Oprea TI** (2004) Pursuing the leadlikeness concept in pharmaceutical  
751 research. *Current Opinion in Chemical Biology* **8**(3), 255-263.  
752 <https://doi.org/10.1016/j.cbpa.2004.04.003>.
- 753 **Hindorff LA, Bonham VL, Brody LC, Ginoza MEC, Hutter CM, Manolio TA and Green**  
754 **ED** (2018) Prioritizing diversity in human genomics research. *Nature Reviews*  
755 *Genetics* **19**(3), 175-185. <https://doi.org/10.1038/nrg.2017.89>.
- 756 **Ho D, Yan L, Iwatsubo K, Vatner DE and Vatner SF** (2010) Modulation of  $\beta$ -  
757 adrenergic receptor signaling in heart failure and longevity: targeting adenylyl  
758 cyclase type 5. *Heart Failure Reviews* **15**(5), 495-512.  
759 <https://doi.org/10.1007/s10741-010-9183-5>.
- 760 **Huang J, Zhu H, Haggarty SJ, Spring DR, Hwang H, Jin F, Snyder M and Schreiber**  
761 **SL** (2004) Finding new components of the target of rapamycin (TOR) signaling

- 762 network through chemical genetics and proteome chips. *Proceedings of the*  
763 *National Academy of Sciences* **101**(47), 16594-16599.  
764 <https://doi.org/10.1073/pnas.0407117101>.
- 765 **Huang K, Fu T, Glass LM, Zitnik M, Xiao C and Sun J (2021)** DeepPurpose: a deep  
766 learning library for drug-target interaction prediction. *Bioinformatics* **36**(22-23),  
767 5545-5547. <https://doi.org/10.1093/bioinformatics/btaa1005>.
- 768 **Ingraham JB, Baranov M, Costello Z, Barber KW, Wang W, Ismail A, Frappier V,**  
769 **Lord DM, Ng-Thow-Hing C, Van Vlack ER, Tie S, Xue V, Cowles SC, Leung**  
770 **A, Rodrigues JV, Morales-Perez CL, Ayoub AM, Green R, Puentes K,**  
771 **Oplinger F, Panwar NV, Obermeyer F, Root AR, Beam AL, Poelwijk FJ and**  
772 **Grigoryan G (2023)** Illuminating protein space with a programmable  
773 generative model. *Nature* 1-9. <https://doi.org/10.1038/s41586-023-06728-8>.
- 774 **Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O,**  
775 **Tunyasuvunakool K, Bates R, Židek A, Potapenko A, Bridgland A, Meyer C,**  
776 **Kohl SAA, Ballard AJ, Cowie A, Romera-Paredes B, Nikolov S, Jain R, Adler**  
777 **J, Back T, Petersen S, Reiman D, Clancy E, Zielinski M, Steinegger M,**  
778 **Pacholska M, Berghammer T, Bodenstein S, Silver D, Vinyals O, Senior AW,**  
779 **Kavukcuoglu K, Kohli P and Hassabis D (2021)** Highly accurate protein  
780 structure prediction with AlphaFold. *Nature* **596**(7873), 583-589.  
781 <https://doi.org/10.1038/s41586-021-03819-2>.
- 782 **Kanehisa M, Furumichi M, Sato Y, Kawashima M and Ishiguro-Watanabe M (2023)**  
783 KEGG for taxonomy-based analysis of pathways and genomes. *Nucleic Acids*  
784 *Research* **51**(D1), D587-D592. <https://doi.org/10.1093/nar/gkac963>.



- 785 **Kedzierska KZ, Crawford L, Amini AP and Lu AX** (2023, November 5) Assessing the  
786 limits of zero-shot foundation models in single-cell biology. *bioRxiv*,  
787 2023.10.16.561085. <https://doi.org/10.1101/2023.10.16.561085>.
- 788 **Kelley DR, Reshef Y, Bileschi M, Belanger D, McLean CY and Snoek J** (2018)  
789 Sequential regulatory activity prediction across chromosomes with  
790 convolutional neural networks. *Genome Research* gr.227819.117.  
791 <https://doi.org/10.1101/gr.227819.117>.
- 792 **Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM and Haussler**  
793 **and D** (2002) The Human Genome Browser at UCSC. *Genome Research*  
794 12(6), 996-1006. <https://doi.org/10.1101/gr.229102>.
- 795 **Kerrien S, Alam-Faruque Y, Aranda B, Bancarz I, Bridge A, Derow C, Dimmer E,**  
796 **Feuermann M, Friedrichsen A, Huntley R, Kohler C, Khadake J, Leroy C,**  
797 **Liban A, Liefink C, Montecchi-Palazzi L, Orchard S, Risse J, Robbe K,**  
798 **Roechert B, Thorneycroft D, Zhang Y, Apweiler R and Hermjakob H** (2007)  
799 IntAct—open source resource for molecular interaction data. *Nucleic Acids*  
800 *Research* 35(suppl\_1), D561-D565. <https://doi.org/10.1093/nar/gkl958>.
- 801 **Khader F, Kather JN, Müller-Franzes G, Wang T, Han T, Tayebi Arasteh S,**  
802 **Hamesch K, Bressemer K, Haarbuerger C, Stegmaier J, Kuhl C, Nebelung S and**  
803 **Truhn D** (2023) Medical transformer for multimodal survival prediction in  
804 intensive care: integration of imaging and non-imaging data. *Scientific Reports*  
805 13(1), 10666. <https://doi.org/10.1038/s41598-023-37835-1>.
- 806 **Krentzel D, Shorte SL and Zimmer C** (2023) Deep learning in image-based  
807 phenotypic drug discovery. *Trends in Cell Biology* 33(7), 538-554.  
808 <https://doi.org/10.1016/j.tcb.2022.11.011>.

- 809 **Kuhn M, von Mering C, Campillos M, Jensen LJ and Bork P (2008)** STITCH:  
810 interaction networks of chemicals and proteins. *Nucleic Acids Research*  
811 **36**(Database issue), D684-D688. <https://doi.org/10.1093/nar/gkm795>.
- 812 **Kursa MB, Jankowski A and Rudnicki WR (2010)** Boruta - A System for Feature  
813 Selection. *Fundamenta Informaticae* **101**(4), 271-285.  
814 <https://doi.org/10.3233/FI-2010-288>.
- 815 **Landry LG, Ali N, Williams DR, Rehm HL and Bonham VL (2018)** Lack Of Diversity In  
816 Genomic Databases Is A Barrier To Translating Precision Medicine Research  
817 Into Practice. *Health Affairs* **37**(5), 780-785.  
818 <https://doi.org/10.1377/hlthaff.2017.1595>.
- 819 **LeCun Y, Bengio Y and Hinton G (2015)** Deep learning. *Nature* **521**(7553), 436-444.  
820 <https://doi.org/10.1038/nature14539>.
- 821 **Leiserson MDM, Vandin F, Wu H-T, Dobson JR, Eldridge JV, Thomas JL,**  
822 **Papoutsaki A, Kim Y, Niu B, McLellan M, Lawrence MS, Gonzalez-Perez A,**  
823 **Tamborero D, Cheng Y, Ryslik GA, Lopez-Bigas N, Getz G, Ding L and**  
824 **Raphael BJ (2015)** Pan-cancer network analysis identifies combinations of  
825 rare somatic mutations across pathways and protein complexes. *Nature*  
826 *Genetics* **47**(2), 106-114. <https://doi.org/10.1038/ng.3168>.
- 827 **Lin S, Shi C and Chen J (2022)** GeneralizedDTA: combining pre-training and multi-  
828 task learning to predict drug-target binding affinity for unknown drug  
829 discovery. *BMC Bioinformatics* **23**(1), 367. [https://doi.org/10.1186/s12859-](https://doi.org/10.1186/s12859-022-04905-6)  
830 [022-04905-6](https://doi.org/10.1186/s12859-022-04905-6).
- 831 **Lin W, Wells J, Wang Z, Orengo C and Martin ACR (2023a, March 20)** VariPred:  
832 Enhancing Pathogenicity Prediction of Missense Variants Using Protein

- 833 Language Models. bioRxiv, 2023.03.16.532942.  
834 <https://doi.org/10.1101/2023.03.16.532942>.
- 835 **Lin Z, Akin H, Rao R, Hie B, Zhu Z, Lu W, Smetanin N, Verkuil R, Kabeli O, Shmueli**  
836 **Y, dos Santos Costa A, Fazel-Zarandi M, Sercu T, Candido S and Rives A**  
837 **(2023b) Evolutionary-scale prediction of atomic-level protein structure with a**  
838 **language model. *Science* 379(6637), 1123-1130.**  
839 <https://doi.org/10.1126/science.ade2574>.
- 840 **Linder J, Srivastava D, Yuan H, Agarwal V and Kelley DR (2023, September 1)**  
841 **Predicting RNA-seq coverage from DNA sequence as a unifying model of**  
842 **gene regulation. bioRxiv, 2023.08.30.555582.**  
843 <https://doi.org/10.1101/2023.08.30.555582>.
- 844 **Liu Y, Chen X, Cheng J and Peng H (2017) A medical image fusion method based**  
845 **on convolutional neural networks. In *2017 20th International Conference on***  
846 ***Information Fusion (Fusion)*. 1-7. <https://doi.org/10.23919/ICIF.2017.8009769>.**
- 847 **Liu Z-P, Wu C, Miao H and Wu H (2015) RegNetwork: an integrated database of**  
848 **transcriptional and post-transcriptional regulatory networks in human and**  
849 **mouse. *Database* 2015, bav095. <https://doi.org/10.1093/database/bav095>.**
- 850 **Loynachan C, Unsworth H, Donoghue K and Sonabend R (2023) Unlocking the**  
851 **Potential of AI in Drug Discovery. *Boston Consulting Group and Wellcome*.**
- 852 **Martin FJ, Amode MR, Aneja A, Austine-Orimoloye O, Azov AG, Barnes I, Becker A,**  
853 **Bennett R, Berry A, Bhai J, Bhurji SK, Bignell A, Boddu S, Branco Lins PR,**  
854 **Brooks L, Ramaraju SB, Charkhchi M, Cockburn A, Da Rin Fiorretto L,**  
855 **Davidson C, Dodiya K, Donaldson S, El Houdaigui B, El Naboulsi T, Fatima R,**  
856 **Giron CG, Genez T, Ghattaoraya GS, Martinez JG, Guijarro C, Hardy M,**  
857 **Hollis Z, Hourlier T, Hunt T, Kay M, Kaykala V, Le T, Lemos D, Marques-**

858 Coelho D, Marugán JC, Merino GA, Mirabueno LP, Mushtaq A, Hossain SN,  
859 Ogeh DN, Sakthivel MP, Parker A, Perry M, Piližota I, Prosovetskaia I, Pérez-  
860 Silva JG, Salam AIA, Saraiva-Agostinho N, Schuilenburg H, Sheppard D,  
861 Sinha S, Sipos B, Stark W, Steed E, Sukumaran R, Sumathipala D, Suner M-  
862 M, Surapaneni L, Sutinen K, Szpak M, Tricomi FF, Urbina-Gómez D,  
863 Veidenberg A, Walsh TA, Walts B, Wass E, Willhoft N, Allen J, Alvarez-  
864 Jarreta J, Chakiachvili M, Flint B, Giorgetti S, Haggerty L, Ilsley GR, Loveland  
865 JE, Moore B, Mudge JM, Tate J, Thybert D, Trevanion SJ, Winterbottom A,  
866 Frankish A, Hunt SE, Ruffier M, Cunningham F, Dyer S, Finn RD, Howe KL,  
867 Harrison PW, Yates AD and Flicek P (2023) Ensembl 2023. *Nucleic Acids  
868 Research* 51(D1), D933-D941. <https://doi.org/10.1093/nar/gkac958>.

869 **Materi W and Wishart DS (2007)** Computational systems biology in drug discovery  
870 and development: methods and applications. *Drug Discovery Today* 12(7),  
871 295-303. <https://doi.org/10.1016/j.drudis.2007.02.013>.

872 **May LT, Leach K, Sexton PM and Christopoulos A (2007)** Allosteric Modulation of G  
873 Protein-Coupled Receptors. *Annual Review of Pharmacology and Toxicology*  
874 47(1), 1-51. <https://doi.org/10.1146/annurev.pharmtox.47.120505.105159>.

875 **Mendez D, Gaulton A, Bento AP, Chambers J, De Veij M, Félix E, Magariños MP,  
876 Mosquera JF, Mutowo P, Nowotka M, Gordillo-Marañón M, Hunter F, Junco L,  
877 Mugumbate G, Rodriguez-Lopez M, Atkinson F, Bosc N, Radoux CJ, Segura-  
878 Cabrera A, Hersey A and Leach AR (2019)** ChEMBL: towards direct  
879 deposition of bioassay data. *Nucleic Acids Research* 47(D1), D930-D940.  
880 <https://doi.org/10.1093/nar/gky1075>.

- 881 **Microsoft Research AI4Science and Microsoft Azure Quantum** (2023, November 13)  
882 The Impact of Large Language Models on Scientific Discovery: a Preliminary  
883 Study using GPT-4. arXiv. <https://doi.org/10.48550/arXiv.2311.07361>.
- 884 **Mikolov T, Chen K, Corrado G and Dean J** (2013a, September 6) Efficient Estimation  
885 of Word Representations in Vector Space. arXiv.  
886 <https://doi.org/10.48550/arXiv.1301.3781>.
- 887 **Mikolov T, Sutskever I, Chen K, Corrado GS and Dean J** (2013b) Distributed  
888 Representations of Words and Phrases and their Compositionality. In  
889 *Advances in Neural Information Processing Systems*, Vol. 26. Curran  
890 Associates, Inc.  
891 [https://proceedings.neurips.cc/paper\\_files/paper/2013/hash/9aa42b31882ec0](https://proceedings.neurips.cc/paper_files/paper/2013/hash/9aa42b31882ec039965f3c4923ce901b-Abstract.html)  
892 [39965f3c4923ce901b-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2013/hash/9aa42b31882ec039965f3c4923ce901b-Abstract.html) (accessed 18 September 2023)
- 893 **Minikel EV, Painter JL, Dong CC and Nelson MR** (2024) Refining the impact of  
894 genetic evidence on clinical success. *Nature* 1-6.  
895 <https://doi.org/10.1038/s41586-024-07316-0>.
- 896 **Moffat JG, Vincent F, Lee JA, Eder J and Prunotto M** (2017) Opportunities and  
897 challenges in phenotypic drug discovery: an industry perspective. *Nature*  
898 *Reviews Drug Discovery* **16**(8), 531-543. <https://doi.org/10.1038/nrd.2017.111>.
- 899 **Muzio G, O'Bray L and Borgwardt K** (2021) Biological network analysis with deep  
900 learning. *Briefings in Bioinformatics* **22**(2), 1515-1530.  
901 <https://doi.org/10.1093/bib/bbaa257>.
- 902 **Nakagawa O, Arnold M, Nakagawa M, Hamada H, Shelton JM, Kusano H, Harris**  
903 **TM, Childs G, Campbell KP, Richardson JA, Nishino I and Olson EN** (2005)  
904 Centronuclear myopathy in mice lacking a novel muscle-specific protein

905 kinase transcriptionally regulated by MEF2. *Genes & Development* **19**(17),  
906 2066-2077. <https://doi.org/10.1101/gad.1338705>.

907 **Ngiam J, Khosla A, Kim M, Nam J, Lee H and Ng AY** (2011) Multimodal deep  
908 learning. In *Proceedings of the 28th International Conference on International*  
909 *Conference on Machine Learning*. Madison, WI, USA: Omnipress, 689-696.  
910 (accessed 23 November 2023)

911 **Ochoa D, Hercules A, Carmona M, Suveges D, Gonzalez-Uriarte A, Malangone C,**  
912 **Miranda A, Fumis L, Carvalho-Silva D, Spitzer M, Baker J, Ferrer J, Raies A,**  
913 **Razuvayevskaya O, Faulconbridge A, Petsalaki E, Mutowo P, Machlitt-**  
914 **Northen S, Peat G, McAuley E, Ong CK, Mountjoy E, Ghossaini M, Pierleoni**  
915 **A, Papa E, Pignatelli M, Koscielny G, Karim M, Schwartzentruber J, Hulcoop**  
916 **DG, Dunham I and McDonagh EM** (2021) Open Targets Platform: supporting  
917 systematic drug-target identification and prioritisation. *Nucleic Acids Research*  
918 **49**(D1), D1302-D1310. <https://doi.org/10.1093/nar/gkaa1027>.

919 **Oprea TI, Bologa CG, Brunak S, Campbell A, Gan GN, Gaulton A, Gomez SM, Guha**  
920 **R, Hersey A, Holmes J, Jadhav A, Jensen LJ, Johnson GL, Karlson A, Leach**  
921 **AR, Ma'ayan A, Malovannaya A, Mani S, Mathias SL, McManus MT, Meehan**  
922 **TF, von Mering C, Muthas D, Nguyen D-T, Overington JP, Papadatos G, Qin**  
923 **J, Reich C, Roth BL, Schürer SC, Simeonov A, Sklar LA, Southall N, Tomita**  
924 **S, Tudose I, Ursu O, Vidović D, Waller A, Westergaard D, Yang JJ and**  
925 **Zahoránszky-Köhalmi G** (2018) Unexplored therapeutic opportunities in the  
926 human genome. *Nature Reviews Drug Discovery* **17**(5), 317-332.  
927 <https://doi.org/10.1038/nrd.2018.14>.

- 928 **Paananen J and Fortino V** (2020) An omics perspective on drug target discovery  
929 platforms. *Briefings in Bioinformatics* **21**(6), 1937-1953.  
930 <https://doi.org/10.1093/bib/bbz122>.
- 931 **Pan SJ and Yang Q** (2010) A Survey on Transfer Learning. *IEEE Transactions on*  
932 *Knowledge and Data Engineering* **22**(10), 1345-1359.  
933 <https://doi.org/10.1109/TKDE.2009.191>.
- 934 **Paul SM, Mytelka DS, Dunwiddie CT, Persinger CC, Munos BH, Lindborg SR and**  
935 **Schacht AL** (2010) How to improve R&D productivity: the pharmaceutical  
936 industry's grand challenge. *Nature Reviews. Drug Discovery* **9**(3), 203-214.  
937 <https://doi.org/10.1038/nrd3078>.
- 938 **Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M,**  
939 **Prettenhofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D,**  
940 **Brucher M, Perrot M and Duchesnay É** (2011) Scikit-learn: Machine Learning  
941 in Python. *The Journal of Machine Learning Research* **12**(null), 2825-2830.
- 942 **Polikar R** (2006) Ensemble based systems in decision making. *IEEE Circuits and*  
943 *Systems Magazine* **6**(3), 21-45. <https://doi.org/10.1109/MCAS.2006.1688199>.
- 944 **Popejoy AB and Fullerton SM** (2016) Genomics is failing on diversity. *Nature*  
945 **538**(7624), 161-164. <https://doi.org/10.1038/538161a>.
- 946 **Prasad N, Yang K and Uhler C** (2022, January 5) Optimal Transport using GANs for  
947 Lineage Tracing. arXiv. <https://doi.org/10.48550/arXiv.2007.12098>.
- 948 **Radford A, Narasimhan K, Salimans T and Sutskever I** (2018) Improving Language  
949 Understanding by Generative Pre-Training.
- 950 **Raies A, Tulodziecka E, Stainer J, Middleton L, Dhindsa RS, Hill P, Engkvist O,**  
951 **Harper AR, Petrovski S and Vitsios D** (2022) DrugnomeAI is an ensemble  
952 machine-learning framework for predicting druggability of candidate drug



- 953 targets. *Communications Biology* 5(1), 1-16. <https://doi.org/10.1038/s42003->  
954 022-04245-4.
- 955 **Ramamoorthy A, Pacanowski M, Bull J and Zhang L** (2015) Racial/ethnic differences  
956 in drug disposition and response: Review of recently approved drugs. *Clinical*  
957 *Pharmacology & Therapeutics* 97(3), 263-273. <https://doi.org/10.1002/cpt.61>.
- 958 **Razuvayevskaya O, Lopez I, Dunham I and Ochoa D** (2023, February 8) Why  
959 Clinical Trials Stop: The Role of Genetics. medRxiv, 2023.02.07.23285407.  
960 <https://doi.org/10.1101/2023.02.07.23285407>.
- 961 **Reymond J-L** (2015) The Chemical Space Project. *Accounts of Chemical Research*  
962 48(3), 722-730. <https://doi.org/10.1021/ar500432k>.
- 963 **Reyna MA, Leiserson MDM and Raphael BJ** (2018) Hierarchical HotNet: identifying  
964 hierarchies of altered subnetworks. *Bioinformatics* 34(17), i972-i980.  
965 <https://doi.org/10.1093/bioinformatics/bty613>.
- 966 **Rocha JJ, Jayaram SA, Stevens TJ, Muschalik N, Shah RD, Emran S, Robles C,**  
967 **Freeman M and Munro S** (2023) Functional unknowns: Systematic screening  
968 of conserved genes of unknown function. *PLOS Biology* 21(8), e3002222.  
969 <https://doi.org/10.1371/journal.pbio.3002222>.
- 970 **Royer J, Rodríguez-Cruces R, Tavakol S, Larivière S, Herholz P, Li Q, Vos de Wael**  
971 **R, Paquola C, Benkarim O, Park B, Lowe AJ, Margulies D, Smallwood J,**  
972 **Bernasconi A, Bernasconi N, Frauscher B and Bernhardt BC** (2022) An Open  
973 MRI Dataset For Multiscale Neuroscience. *Scientific Data* 9(1), 569.  
974 <https://doi.org/10.1038/s41597-022-01682-y>.
- 975 **Ruddigkeit L, van Deursen R, Blum LC and Reymond J-L** (2012) Enumeration of 166  
976 Billion Organic Small Molecules in the Chemical Universe Database GDB-17.

- 977 *Journal of Chemical Information and Modeling* **52**(11), 2864-2875.  
978 <https://doi.org/10.1021/ci300415d>.
- 979 **Sadawi N, Olier I, Vanschoren J, van Rijn JN, Besnard J, Bickerton R, Grosan C,**  
980 **Soldatova L and King RD** (2019) Multi-task learning with a natural metric for  
981 quantitative structure activity relationship learning. *Journal of Cheminformatics*  
982 **11**(1), 68. <https://doi.org/10.1186/s13321-019-0392-1>.
- 983 **Santos R, Ursu O, Gaulton A, Bento AP, Donadi RS, Bologa CG, Karlsson A, Al-**  
984 **Lazikani B, Hersey A, Oprea TI and Overington JP** (2017) A comprehensive  
985 map of molecular drug targets. *Nature Reviews. Drug Discovery* **16**(1), 19-34.  
986 <https://doi.org/10.1038/nrd.2016.230>.
- 987 **Sayers EW, Cavanaugh M, Clark K, Pruitt KD, Schoch CL, Sherry ST and Karsch-**  
988 **Mizrachi I** (2021) GenBank. *Nucleic Acids Research* **49**(D1), D92-D96.  
989 <https://doi.org/10.1093/nar/gkaa1023>.
- 990 **Scannell JW, Blanckley A, Boldon H and Warrington B** (2012) Diagnosing the decline  
991 in pharmaceutical R&D efficiency. *Nature Reviews. Drug Discovery* **11**(3),  
992 191-200. <https://doi.org/10.1038/nrd3681>.
- 993 **Schulte-Sasse R, Budach S, Hnisz D and Marsico A** (2021) Integration of multiomics  
994 data with graph convolutional networks to identify new cancer genes and their  
995 associated molecular mechanisms. *Nature Machine Intelligence* **3**(6), 513-  
996 526. <https://doi.org/10.1038/s42256-021-00325-y>.
- 997 **Shabani M and Hojjat-Farsangi M** (2016) Targeting Receptor Tyrosine Kinases  
998 Using Monoclonal Antibodies: The Most Specific Tools for Targeted-Based  
999 Cancer Therapy. *Current Drug Targets* **17**(14), 1687-1703.
- 1000 **Singhal K, Azizi S, Tu T, Mahdavi SS, Wei J, Chung HW, Scales N, Tanwani A,**  
1001 **Cole-Lewis H, Pfohl S, Payne P, Seneviratne M, Gamble P, Kelly C, Babiker**

1002 A, Schärli N, Chowdhery A, Mansfield P, Demner-Fushman D, Agüera y Arcas  
1003 B, Webster D, Corrado GS, Matias Y, Chou K, Gottweis J, Tomasev N, Liu Y,  
1004 Rajkomar A, Barral J, Semturs C, Karthikesalingam A and Natarajan V  
1005 (2023a) Large language models encode clinical knowledge. *Nature*  
1006 620(7972), 172-180. <https://doi.org/10.1038/s41586-023-06291-2>.

1007 Singhal K, Tu T, Gottweis J, Sayres R, Wulczyn E, Hou L, Clark K, Pfohl S, Cole-  
1008 Lewis H, Neal D, Schaekermann M, Wang A, Amin M, Lachgar S, Mansfield  
1009 P, Prakash S, Green B, Dominowska E, Arcas BA y, Tomasev N, Liu Y, Wong  
1010 R, Semturs C, Mahdavi SS, Barral J, Webster D, Corrado GS, Matias Y, Azizi  
1011 S, Karthikesalingam A and Natarajan V (2023b, May 16) Towards Expert-  
1012 Level Medical Question Answering with Large Language Models. arXiv.  
1013 <https://doi.org/10.48550/arXiv.2305.09617>.

1014 Stark C, Breitkreutz B-J, Reguly T, Boucher L, Breitkreutz A and Tyers M (2006)  
1015 BioGRID: a general repository for interaction datasets. *Nucleic Acids*  
1016 *Research* 34(suppl\_1), D535-D539. <https://doi.org/10.1093/nar/gkj109>.

1017 Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, Downey P, Elliott P,  
1018 Green J, Landray M, Liu B, Matthews P, Ong G, Pell J, Silman A, Young A,  
1019 Sprosen T, Peakman T and Collins R (2015) UK Biobank: An Open Access  
1020 Resource for Identifying the Causes of a Wide Range of Complex Diseases of  
1021 Middle and Old Age. *PLOS Medicine* 12(3), e1001779.  
1022 <https://doi.org/10.1371/journal.pmed.1001779>.

1023 Szklarczyk D, Kirsch R, Koutrouli M, Nastou K, Mehryary F, Hachilif R, Gable AL,  
1024 Fang T, Doncheva NT, Pyysalo S, Bork P, Jensen LJ and von Mering C  
1025 (2023) The STRING database in 2023: protein-protein association networks  
1026 and functional enrichment analyses for any sequenced genome of interest.

- 1027 *Nucleic Acids Research* **51**(D1), D638-D646.  
1028 <https://doi.org/10.1093/nar/gkac1000>.
- 1029 **Terstappen GC, Schlüpen C, Raggiaschi R and Gaviraghi G (2007)** Target  
1030 deconvolution strategies in drug discovery. *Nature Reviews Drug Discovery*  
1031 **6**(11), 891-903. <https://doi.org/10.1038/nrd2410>.
- 1032 **The UniProt Consortium (2023)** UniProt: the Universal Protein Knowledgebase in  
1033 2023. *Nucleic Acids Research* **51**(D1), D523-D531.  
1034 <https://doi.org/10.1093/nar/gkac1052>.
- 1035 **Theodoris CV, Xiao L, Chopra A, Chaffin MD, Al Sayed ZR, Hill MC, Mantineo H,**  
1036 **Brydon EM, Zeng Z, Liu XS and Ellinor PT (2023)** Transfer learning enables  
1037 predictions in network biology. *Nature* **618**(7965), 616-624.  
1038 <https://doi.org/10.1038/s41586-023-06139-9>.
- 1039 **Thomas M, Bender A and de Graaf C (2023)** Integrating structure-based approaches  
1040 in generative molecular design. *Current Opinion in Structural Biology* **79**,  
1041 102559. <https://doi.org/10.1016/j.sbi.2023.102559>.
- 1042 **Tu C, Du D, Zeng T and Zhang Y (2023)** Deep Multi-dictionary Learning for Survival  
1043 Prediction with Multi-zoom Histopathological Whole Slide Images. *IEEE/ACM*  
1044 *Transactions on Computational Biology and Bioinformatics* 1-12.  
1045 <https://doi.org/10.1109/TCBB.2023.3321593>.
- 1046 **Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Ł and**  
1047 **Polosukhin I (2017)** Attention is All you Need. In *Advances in Neural*  
1048 *Information Processing Systems*, Vol. 30. Curran Associates, Inc.  
1049 [https://proceedings.neurips.cc/paper\\_files/paper/2017/hash/3f5ee243547dee9](https://proceedings.neurips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html)  
1050 [1fbd053c1c4a845aa-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html) (accessed 4 September 2023)

- 1051 **Veličković P** (2023) Everything is connected: Graph neural networks. *Current*  
1052 *Opinion in Structural Biology* **79**, 102538.  
1053 <https://doi.org/10.1016/j.sbi.2023.102538>.
- 1054 **Venugopalan J, Tong L, Hassanzadeh HR and Wang MD** (2021) Multimodal deep  
1055 learning models for early detection of Alzheimer's disease stage. *Scientific*  
1056 *Reports* **11**(1), 3254. <https://doi.org/10.1038/s41598-020-74399-w>.
- 1057 **Vitsios D and Petrovski S** (2020) Mantis-ml: Disease-Agnostic Gene Prioritization  
1058 from High-Throughput Genomic Screens by Stochastic Semi-supervised  
1059 Learning. *The American Journal of Human Genetics* **106**(5), 659-678.  
1060 <https://doi.org/10.1016/j.ajhg.2020.03.012>.
- 1061 **Wagner AH, Coffman AC, Ainscough BJ, Spies NC, Skidmore ZL, Campbell KM,**  
1062 **Krysiak K, Pan D, McMichael JF, Eldred JM, Walker JR, Wilson RK, Mardis**  
1063 **ER, Griffith M and Griffith OL** (2016a) DGIdb 2.0: mining clinically relevant  
1064 drug-gene interactions. *Nucleic Acids Research* **44**(D1), D1036-D1044.  
1065 <https://doi.org/10.1093/nar/gkv1165>.
- 1066 **Wagner JR, Lee CT, Durrant JD, Malmstrom RD, Feher VA and Amaro RE** (2016b)  
1067 Emerging Computational Methods for the Rational Discovery of Allosteric  
1068 Drugs. *Chemical Reviews* **116**(11), 6370-6390.  
1069 <https://doi.org/10.1021/acs.chemrev.5b00631>.
- 1070 **Watson JL, Juergens D, Bennett NR, Trippe BL, Yim J, Eisenach HE, Ahern W,**  
1071 **Borst AJ, Ragotte RJ, Milles LF, Wicky BIM, Hanikel N, Pellock SJ, Courbet A,**  
1072 **Sheffler W, Wang J, Venkatesh P, Sappington I, Torres SV, Lauko A, De**  
1073 **Bortoli V, Mathieu E, Ovchinnikov S, Barzilay R, Jaakkola TS, DiMaio F, Baek**  
1074 **M and Baker D** (2023) De novo design of protein structure and function with

- 1075 RFdiffusion. *Nature* **620**(7976), 1089-1100. <https://doi.org/10.1038/s41586->  
1076 023-06415-8.
- 1077 **Weinstein JN, Collisson EA, Mills GB, Shaw KM, Ozenberger BA, Ellrott K,**  
1078 **Shmulevich I, Sander C and Stuart JM** (2013) The Cancer Genome Atlas Pan-  
1079 Cancer Analysis Project. *Nature Genetics* **45**(10), 1113-1120.  
1080 <https://doi.org/10.1038/ng.2764>.
- 1081 **Wen M, Zhang Z, Niu S, Sha H, Yang R, Yun Y and Lu H** (2017) Deep-Learning-  
1082 Based Drug-Target Interaction Prediction. *Journal of Proteome Research*  
1083 **16**(4), 1401-1409. <https://doi.org/10.1021/acs.jproteome.6b00618>.
- 1084 **Wood V, Lock A, Harris MA, Rutherford K, Bähler J and Oliver SG** (2019) Hidden in  
1085 plain sight: what remains to be discovered in the eukaryotic proteome? *Open*  
1086 *Biology* **9**(2), 180241. <https://doi.org/10.1098/rsob.180241>.
- 1087 **Zrimec J, Fu X, Muhammad AS, Skrekas C, Jauniskis V, Speicher NK, Börlin CS,**  
1088 **Verendel V, Chehreghani MH, Dubhashi D, Siewers V, David F, Nielsen J and**  
1089 **Zelezniak A** (2022) Controlling gene expression with deep generative design  
1090 of regulatory DNA. *Nature Communications* **13**(1), 5099.  
1091 <https://doi.org/10.1038/s41467-022-32818-8>.
- 1092

**1093 Impact Statement**

1094 Artificial intelligence (AI) is transforming drug discovery and development by enabling  
1095 the rapid analysis of massive amounts of biological data and chemical information.  
1096 This paper reviews recent advances in using AI methods for the discovery and  
1097 validation of drug targets. Identifying and validating novel drug targets is fundamental  
1098 to creating safe and effective new medicines but has remained a major bottleneck in  
1099 the drug R&D process. By integrating diverse datasets, AI models can accurately  
1100 predict key properties of drug targets, reveal intricate biological relationships  
1101 underlying disease, and guide drug discovery strategies.

1102  
1103 This paper highlights groundbreaking applications of AI that accelerate target  
1104 discovery, including models that prioritise candidate genes, predict druggability of  
1105 proteins, uncover disease mechanisms, and simulate biological experiments. Critically,  
1106 AI enables leveraging insights across modalities like sequences (e.g. DNA, proteins),  
1107 structures (e.g. compounds, proteins), multiomics, biomedical literature and more.  
1108 Integrating multimodal inputs is paramount for comprehensively understanding  
1109 complex diseases involving genetic and non-genetic factors.

1110  
1111 The AI methods outlined will profoundly enhance R&D efficiency. By illuminating novel  
1112 drug targets, AI-powered target discovery will expand treatment options available for  
1113 patients suffering from previously untreatable or poorly managed diseases. From rare  
1114 diseases and refractory cancers to multifactorial neurodegenerative and autoimmune  
1115 conditions, accelerating target discovery through AI has far-reaching therapeutic  
1116 implications. Additionally, safer, more selective drugs developed against AI-predicted  
1117 targets could dramatically improve patient outcomes and quality of life. Overcoming  
1118 existing challenges in AI-based target discovery will be critical to actualising its  
1119 immense potential and promises to usher in a new era of data-driven, accelerated  
1120 drug R&D.



1121 **Financial Support**

1122 C. C. was funded by the National Institute for Health Research (NIHR) as part of the  
1123 portfolio of translational research of the NIHR Biomedical Research Centre at Barts  
1124 and The London School of Medicine and Dentistry. A. W. was funded by the  
1125 UKRI/BBSRC Collaborative Training Partnership in AI for Drug Discovery and Queen  
1126 Mary University of London.

1127

1128 **Conflict of Interest**

1129 At the time of writing, W. W. and V. N. were employed by MSD.

1130

1131

1132

1133

1134