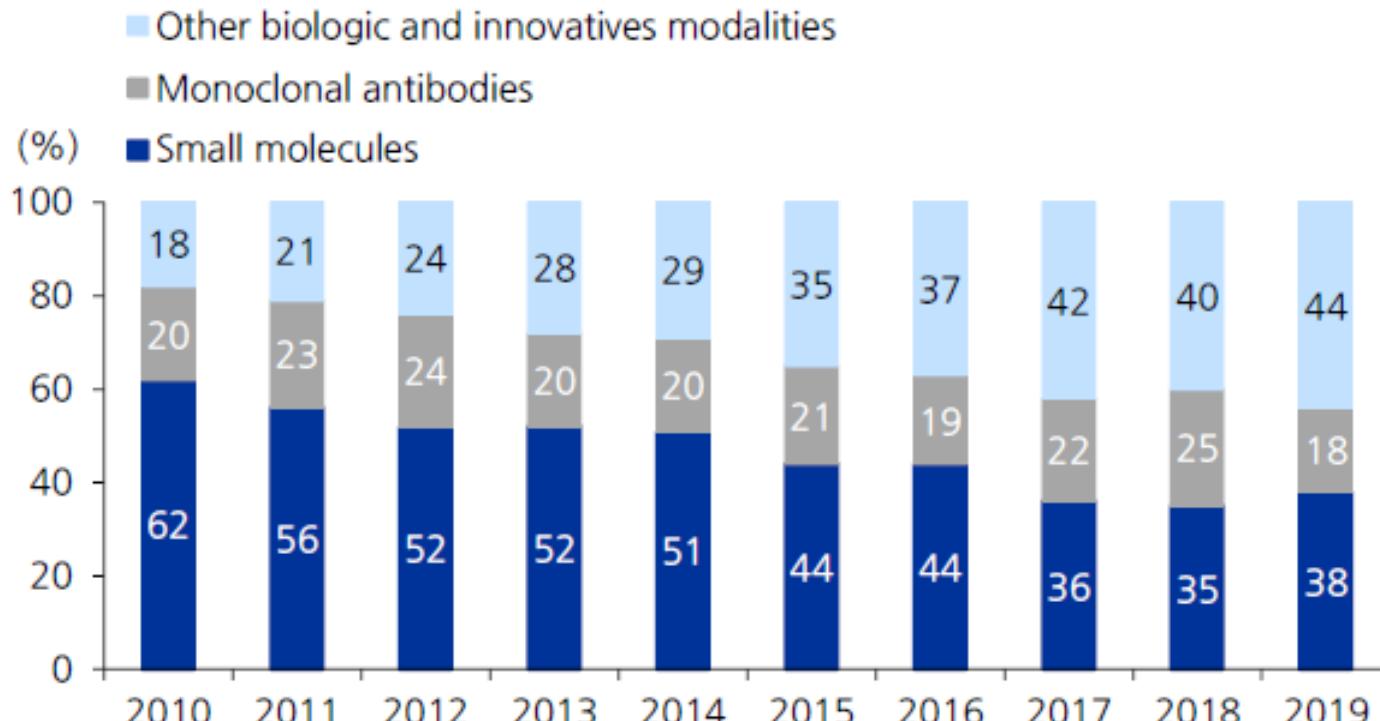


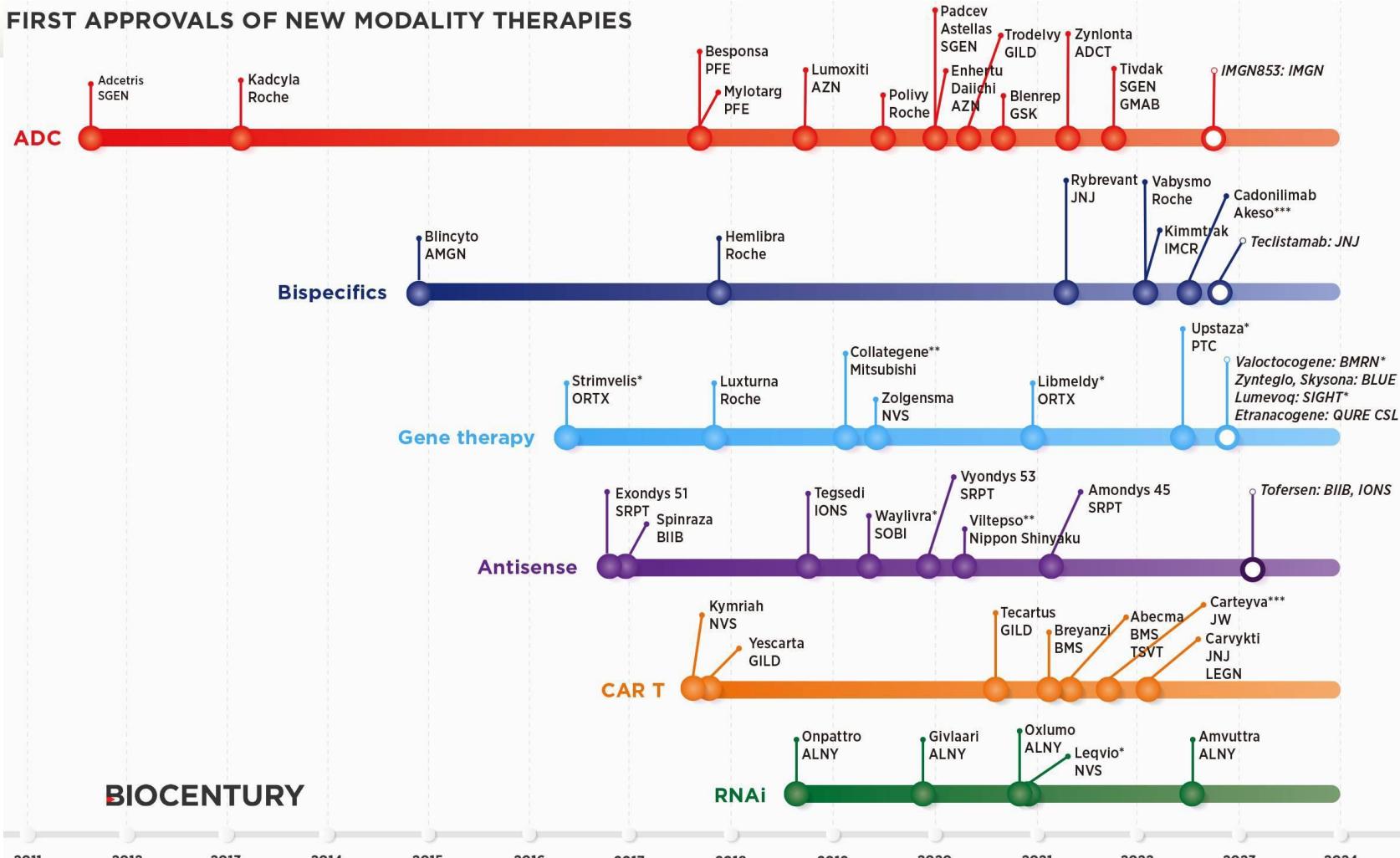
• 07

• 혁신 신약 모달리티의 이해



자료: Reviving an R&D pipelines (2020.10)

## FIRST APPROVALS OF NEW MODALITY THERAPIES



\*Approved/under review in Europe

\*\*Approved in Japan

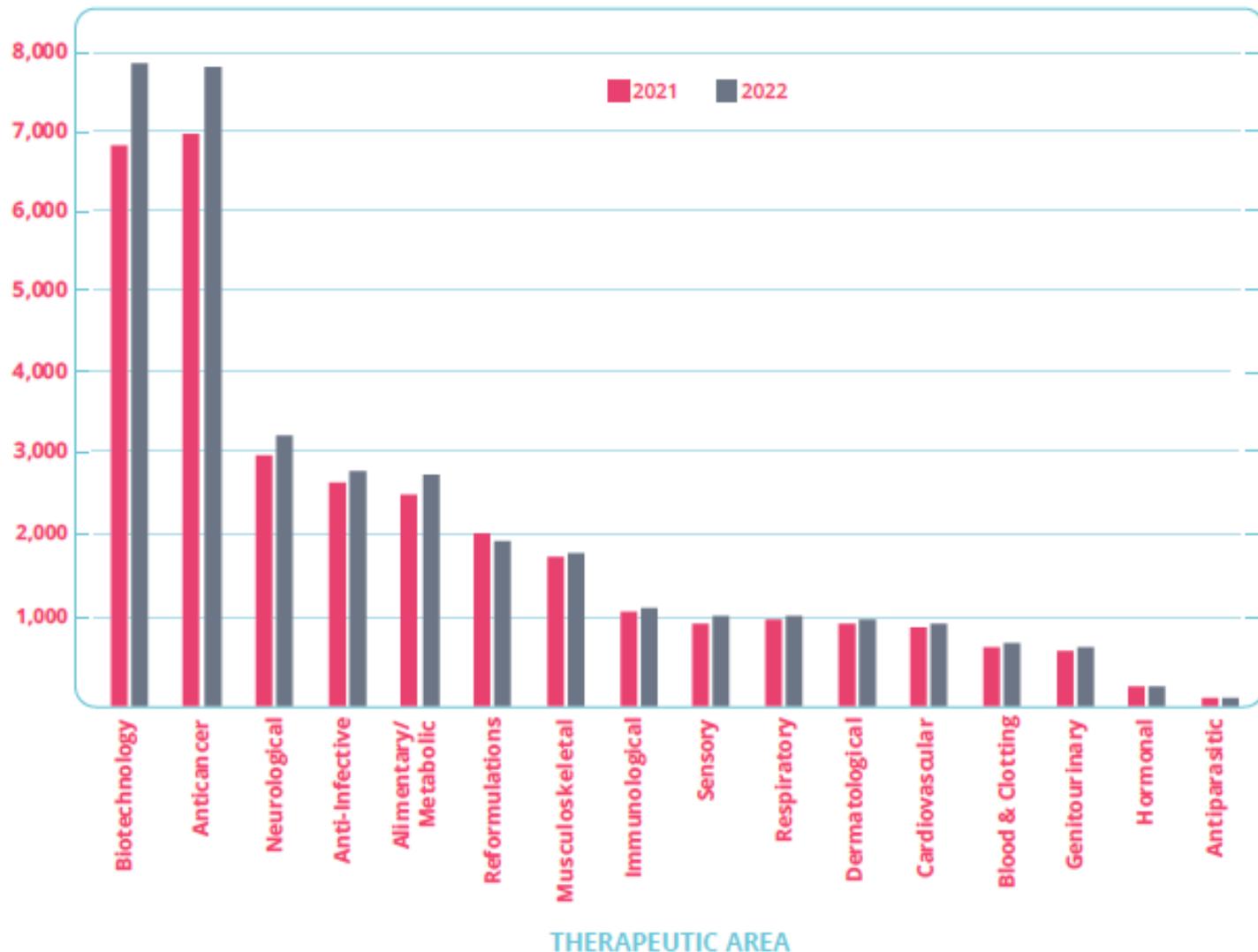
\*\*\*Approved in China

ADC - antibody-drug conjugate

Under regulatory review

Data as of July 31, 2022

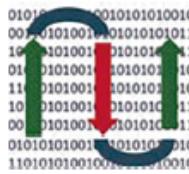
## The R&D pipeline by therapy group, 2021 and 2022



\* Source: xxxx, 2022.0x.xx



ELSEVIER



COMPUTATIONAL  
AND STRUCTURAL  
BIOTECHNOLOGY  
JOURNAL

journal homepage: [www.elsevier.com/locate/csbj](http://www.elsevier.com/locate/csbj)



Review

## Protein–protein interaction prediction with deep learning: A comprehensive review



Farzan Soleymani <sup>a</sup>, Eric Paquet <sup>b,\*</sup>, Herna Viktor <sup>c</sup>, Wojtek Michalowski <sup>d</sup>, Davide Spinello <sup>a</sup>

<sup>a</sup> Department of Mechanical Engineering, University of Ottawa, Ottawa, ON, Canada

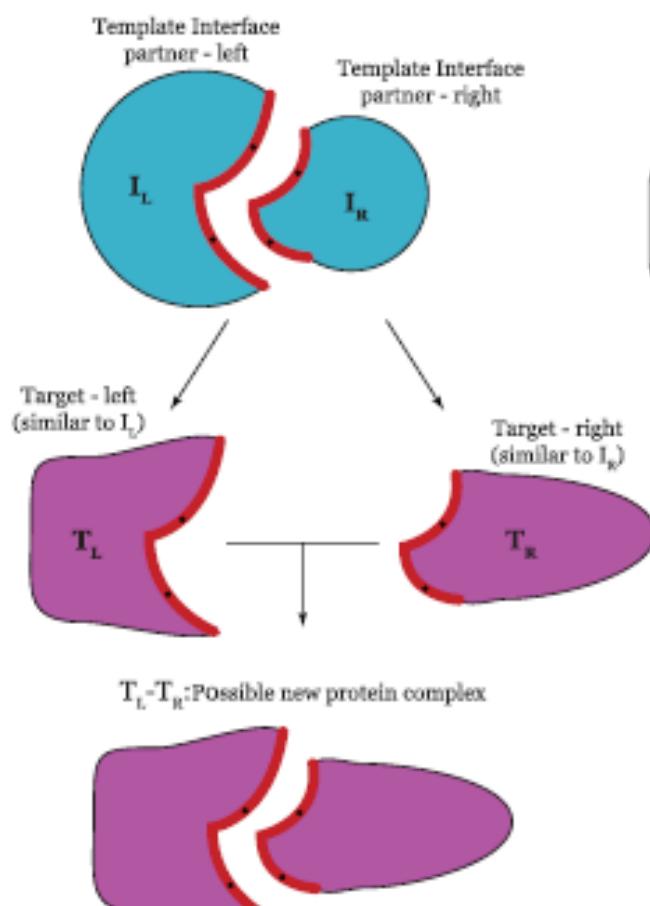
<sup>b</sup> National Research Council, 1200 Montreal Road, Ottawa, ON K1A 0R6, Canada

<sup>c</sup> School of Electrical Engineering and Computer Science, University of Ottawa, ON, Canada

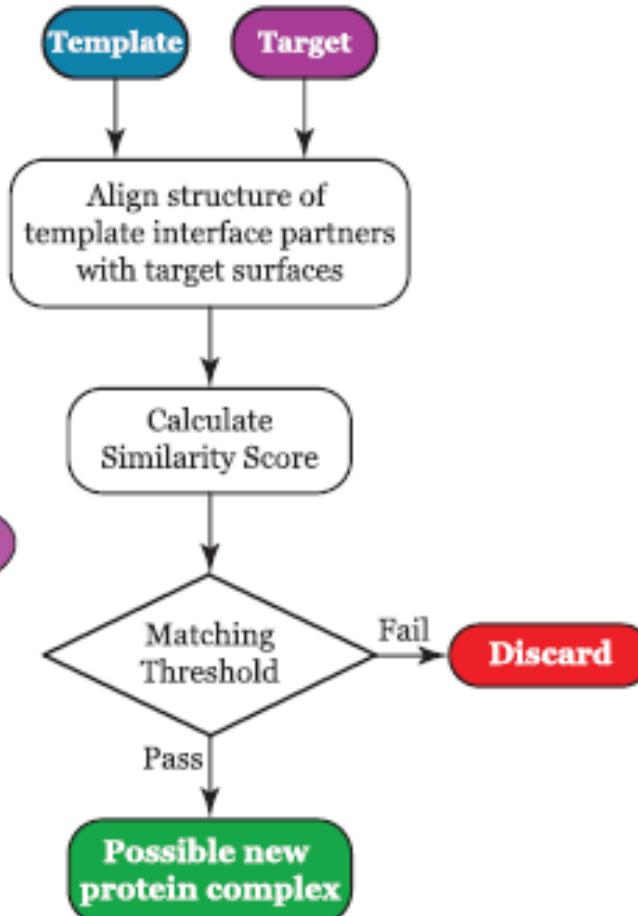
<sup>d</sup> Telfer School of Management, University of Ottawa, ON, K1N 6N5, Canada

**Table 4**  
Databases for PPI prediction

Type	Database	Description	Last update	URL
Protein–Protein Interactions	STRING [109]	Functional associations between protein pairs, which contains 67,592,464 proteins from 14094 organisms; 20,052,394,042 interactions.	2021	<a href="https://string-db.org/">https://string-db.org/</a>
	IntAct [116]	Contains manually curated datasets (topical), interactomes (for 16 different species) and annotations of experimental evidence.	2021	<a href="https://www.ebi.ac.uk/intact/home">https://www.ebi.ac.uk/intact/home</a>
	Biogrid [111]	Contains 2,467,140 protein and genetic interactions, 29,417 chemical interactions and 1,128,339 PTMs from major model organism species.	2020	<a href="http://www.thebiogrid.org/">http://www.thebiogrid.org/</a>
	DIP [112]	Experimentally determined PPI database including biological information of proteins, PPIs and experimental techniques for identifying interactions.	2020	<a href="https://dip.doe-mbi.ucla.edu/dip/Main.cgi">https://dip.doe-mbi.ucla.edu/dip/Main.cgi</a>
	Negatome 2.0 [108]	Contains 21,795 interactions, with scores of zero and one, using text mining from literature and analysing protein complexes from PDB.	2014	<a href="http://mips.helmholtz-muenchen.de/proj/ppi/negatome/">http://mips.helmholtz-muenchen.de/proj/ppi/negatome/</a>
	MINT [114]	Experimentally curated PPI database that includes approximately 117001 PPIs from 607 different species.	2012	<a href="https://mint.bio.uniroma2.it/">https://mint.bio.uniroma2.it/</a>
	HPRD [115]	Consists of 41,327 PPIs, 93,710 PTMs, 22,490 Subcellular Localizations and 112,158 Protein Expressions.	2010	<a href="http://www.hprd.org">http://www.hprd.org</a>
	BIND [113]	PPIs collected from of humans, yeasts, nematodes, etc.	2005	<a href="http://download.baderlab.org/BINDTranslation">http://download.baderlab.org/BINDTranslation</a>
Protein sequences	UniProt [106]	A collection of protein sequence and functional information, including UniProtKB, UniParc, UniRef and Proteomes. UniProtKB contains 567,483 reviewed (Swiss-Prot)—manually annotated, and 231,354,261 unreviewed (TrEMBL)—computationally analysed, protein sequences.	2020	<a href="http://www.uniprot.org">http://www.uniprot.org</a>
	SWISS-MODEL [107]	A web-based integrated service providing information for protein structure homology modelling. The repository contains 2,217,470 models from SWISS-MODEL for UniProtKB targets, as well as 180,107 structures from PDB with mapping to UniProtKB.	2020	<a href="https://swissmodel.expasy.org/">https://swissmodel.expasy.org/</a>
Higher-level structures	PIR [119]	Integrated protein resources, including protein sequences and high-quality annotations by integrating more than 90 biological databases.	2022	<a href="http://pir.georgetown.edu/">http://pir.georgetown.edu/</a>
	RCSB PDB [110]	Information about the 3-D structure of proteins, nucleic acids, and complex assemblies. 191144 structures, 57349 human sequence structures, and 14406 nucleic acid-containing Structures	2021	<a href="https://www.rcsb.org/">https://www.rcsb.org/</a>
	SCOP [122]	Classification of known proteins and a comprehensive description of the structural and evolutionary relationships between them. As of 2022-05-30, this dataset contains 72,448 non-redundant domains, representing 858,316 protein structures.	2022	<a href="http://scop.mrc-lmb.cam.ac.uk/scop">http://scop.mrc-lmb.cam.ac.uk/scop</a>
Genomic information	CGD [124]	A resource for genomic sequence data, genes and protein information for <i>Candida albicans</i> and related species.	2022	<a href="http://www.candidagenome.org/">http://www.candidagenome.org/</a>



(a)



(b)

Framework	Description	Advantage	Disadvantage
GCN-based [2017] [247]	This study proposed a pairwise classification architecture in which one or more graph convolution layers process the neighbourhood of a residue in each protein. Then, the representation of two residues is paired and passed through a dense layer for classification. This study analysed several GCN-based methods, concluding that neighbourhood-based convolution methods outperform diffusion-based convolution and SVM-based methods.	The proposed convolution operators and obtained features may be helpful for other applications, including protein function, catalytic and other functional residues, and protein interactions with DNA and RNA.	The accuracy of this approach is examined based on a limited number of labelled training examples.
IntPred [2018] [240]	This method uses a random forest to predict protein-protein interface sites at both the surface patch and residue levels.	The performance of a binary classifier can be evaluated using different measurements, such as the Matthews' correlation coefficient (MCC), sensitivity, precision, and specificity [240]. IntPred outperformed the methods ProMate [189], PIER [249], PINUP [250], and meta-PPISP [251], but not SPPIDER [252], based on MCC.	The performance of this method depends on the application. For instance, IntPred was better suited to cases in which false positives are less well tolerated than false negatives.
Graph-based generative model [2019] [177]	This method uses a graph transformer model for designing protein sequences given graph representations of 3-D protein structures, leveraging the spatial locality of dependencies in molecular structures.	This method uses a self-attention mechanism to capture higher-order, interaction-based dependencies between sequence and structure. The graph-based model offers computational efficiency due to the representation of long-range sequence dependencies by short-range sequence dependencies in 3-D space [253–255]. Additionally, they achieved linear computational scaling concerning the sequence length and representational flexibility for coarse and fine-grained structure descriptions.	The evaluation dataset only contained chains up to a length of 500, limiting the applicability of this approach.
Struct2Graph [2020] [248]	In this method, graph embeddings of each protein are obtained using an assigned GCN. Next, relevant geometric features associated with query protein pairs are extracted using a mutual attention network. Finally, a feedforward neural network performs a binary classification between interacting and noninteracting pairs.	Struct2Graph only uses 3-D structural information to predict the PPI. They have reported state-of-the-art performance on both balanced and unbalanced datasets.	Limited availability of 3-D structural information may restrict the applicability of this method.
LSTM-based [2020] [187]	The proposed method integrates the 3-D structure and sequence-based information of proteins to predict PPIs. The 3-D coordinate information, hydropathy index, isoelectric point, and amino acid charges of each protein are fed into a pre-trained ResNet50 model to extract features from these attributes. A stacked autoencoder obtains the compact form of encoded proteins using autocovariance and conjoint triad. The structural features from ResNet50 are passed through LSTM and concatenated with features from the stacked autoencoder. The merged features are then fed into the classifier to predict protein pair labels.	This method performs well despite being trained on a low number of instances.	Limited availability of 3-D structural information may restrict the applicability of this method. Additionally, LSTM models are computationally demanding and slow.

Summary of advantages and disadvantages of sequence-based Deep Learning methods regarding PPI prediction.

Framework	Description	Advantage	Disadvantage
SVM-conjoint triad [2007] [286]	<p>Each protein sequence was represented in this study by a vector of amino acid features.</p> <p>The model was developed based on a support vector machine (SVM) integrated with a kernel function and a conjoint triad feature for describing amino acids.</p> <p>This method mapped different types of PPI networks using only sequence information, which could be applied to explore networks for any newly discovered protein with unknown biological relationships.</p> <p>They suggested that methods without local environments for amino acids are often unreliable, so a conjoint triad method was used.</p>	<p>The 20 standard amino acids were clustered into several classes based on their dipoles and side chain volumes to achieve dimensionality reduction of the vector space.</p> <p>This method might predict PPI networks created by pairwise PPIs.</p>	<p>The limited available information on protein pairs restricts the applicability of this method.</p> <p>Additionally, it mainly considers the properties of two nearby amino acids, overlooking long-range interactions.</p>
SVM-autocovariance [2008] [288]	<p>The method combined a new feature representation using autocovariance (AC) and a support vector machine (SVM).</p> <p>AC considers the interactions between more distant amino acids in the protein sequence, specifically long-range interactions.</p> <p>This is an improvement over the method proposed in [286].</p> <p>This model was evaluated using an independent dataset of 11,474 yeast PPIs.</p> <p>This method used a decision tree model, predicting PPIs using only 20 amino acid frequency combinations from interacting and noninteracting proteins as learning features.</p> <p>This study indicated that asparagine, cysteine, and isoleucine frequencies are important features for discerning between interacting and noninteracting protein pairs.</p>	<p>The conjoint triad (CT) method only considered the attributes of an amino acid and its two neighbouring amino acids [286], while long-range interactions are accounted for by the AC method.</p> <p>In this study, AC variables represented information on interactions between one amino acid and its 30 neighbouring amino acids in the protein sequence.</p> <p>This method was scalable due to using a limited number of attributes.</p> <p>Moreover, this method was based on experimentally validated instances from various species, covering many species.</p>	<p>The model achieved a low prediction accuracy of 58.42% in a negative dataset created using the Prcp method [286].</p>
UNISPPPI [2013] [284]	<p>PPI prediction was addressed by integrating a support vector machine (SVM) and a novel matrix-based representation of the sequence order and dipeptide information of the primary protein sequence, extracting more information than amino acid dipeptide composition.</p> <p>The SVM classified the interaction between protein pairs using these feature vectors.</p> <p>The DeepPPI method used a deep neural network architecture network for each protein to extract high-level discriminative features from common protein descriptors.</p> <p>The interaction between two proteins was determined using the one-hot encoding label.</p> <p>This method comprises two different architectures: DeepPPI-Sep, which uses two separate networks as input for each protein, and DeepPPI-Con, which directly links two proteins in a single network.</p>	<p>This method extracts more information hidden in protein primary sequences than amino acid dipeptide composition.</p>	<p>Instances with a classification score of 0.50 were classified as neither PPIs nor non-PPIs,</p> <p>limiting the applicability of this method.</p> <p>Additionally, the obtained accuracy of 79.4% for interacting and 72.6% for noninteracting pairs are relatively low.</p>
SVM-based method [2015] [285]	<p>This method can capture informative features of protein pairs by a layer-wise abstraction.</p> <p>In addition, DeepPPI can automatically learn an internal distributed feature representation from the data.</p>	<p>SVM algorithms performed relatively poorly with noisy data and are unsuitable for large datasets since training time may increase significantly [303].</p> <p>Moreover, finding a proper kernel function was difficult.</p>	<p>The accuracy of DeepPPI for All Human/Yeast dataset are relatively low, and the accuracy of methods proposed in [304] exceeds that of DeepPPI.</p>
Stacked Autoencoder (SAE) [2017] [18]	<p>This method used a stacked autoencoder to predict PPI.</p> <p>The feature extraction from protein sequences was performed using autocovariance (AC) and the conjoint triad (CT).</p>	<p>SAE can learn hidden interaction features of protein sequences.</p>	<p>They used a synthetic negative interaction dataset, and the accuracy of this model for negative interactions is relatively low.</p>
DPII [2018] [298]	<p>This method performed sequence-based PPI prediction using a deep, Siamese-like convolutional neural network combined with random projection and data augmentation.</p> <p>This method captured the composition information, sequential order of amino acids, and co-occurrence of interacting sequence motifs in a protein pair.</p> <p>Each protein was characterised as a probabilistic sequence profile generated by PASI-BLAST.</p> <p>The patterns in each sequence were identified using the convolutional module, comprising multiple layers.</p> <p>The representations learned by the convolutional module were projected to two different spaces using the random projection module, allowing DPPI to explore the combination of protein motifs.</p>	<p>DPII addresses interactions for both homodimeric and heterodimeric proteins.</p> <p>Moreover, this method could model binding affinities.</p>	<p>This method yields lower PPI prediction accuracy on the Scerevisiae core dataset from PIPR based on 5-fold cross-validation compared to PIPR [299] and DeepTrio [302].</p>

Framework	Description	Advantage	Disadvantage
PIPR [2019] [299]	This study proposed an end-to-end framework for PPI prediction based on amino acid sequences using a deep residual recurrent convolutional neural network in the Siamese architecture. This method leveraged an automatic multi-granular feature selection to capture local significant and sequential features from protein sequences.	The Siamese-based learning architecture captured the mutual influence of protein pairs and allowed for generalising to address different PPI prediction tasks without needing predefined features.	RCNN was built using bidirectional gated recurrent units (bidirectional-GRU). However, GRUs suffer from slow convergence and low learning efficiency [305].
S-VGAE [2020] [300]	This model proposed a signed variational graph autoencoder (S-VGAE) that combined sequence information and graph structure. In this method, the PPI network was regarded as an undirected graph. This framework comprised three parts. First, coding the raw protein sequences. Second, the S-VGAE model extracted vector embedding for each protein with sequence information and graph structure. Finally, a simple three-layer softmax classifier. This model was inspired by the variational graph autoencoder (VGAE) [306] that uses latent variables to learn interpretable representations for undirected graphs.	In this method, the cost function was modified only to consider highly confident interactions, making it more robust to noise.	This model used the conjoint triad (CT) method [286] to encode amino acids. However, CT does not account for long-range interactions in the protein sequence.
ACT-SVM [2020] [289]	This method performed feature extraction on the protein sequence to obtain a vector, composition, and transition descriptor and integrated them into a vector. Then, the feature vector was fed into the SVM classifier. The performance of this method was evaluated using 5-fold, 8-fold, and 10-fold cross-validation on H. pylori and human datasets.	They have observed that SVM method outperforms K-Nearest Neighbour (KNN), ANN, RFM, Naive Bayes, Logistic Regression, s for the H. pylori protein pairs.	Finding the proper kernel and hyperparameters was challenging, and training time for SVM classifiers increases with dataset size [307].
D-SCRIPT [2021] [165]	Deep sequence contact residue interaction prediction transfer (D-script) is an interpretable deep learning method generating structurally informative features given protein sequences using a pre-trained language model from [290]. This method used projection modules to reduce the dimension of features, including the residue-contact map of the protein. Finally, the interaction probability was predicted based on the contact maps.	D-SCRIPT generalised to new species considering the sparsity of training data for most model organisms (i.e., it was relatively accurate for cross-species PPI prediction).	Despite its performance for cross-species PPI prediction, D-SCRIPT underperformed on within-species evaluations. The training dataset only included proteins with 50–800 amino acids, limiting the applicability of this method.
SPNet [2021] [9]	The Siamese pyramid network (SPNet) architecture used self-binding and folding amino acid sequences to predict the binding probability for two proteins based solely on their amino acid sequences. Subsequent screening through potential candidates was performed based on binding probabilities.	This architecture consisted of a multilevel pyramid feature structure encompassing various PPI mechanisms to reduce gradient explosion and disappearance, a multilevel Siamese neural network with an attention mechanism, and a multilevel, trainable binding probability prediction network.	LSTM models are computationally demanding and slow. Moreover, a large number of trees in the random forest leads to a longer training time.
BILSTM-RF [2021] [291]	The BILSTM-RF model extracted features of protein pairs in the human database. BILSTM comprises forward and backwards LSTMs and is capable of bidirectional encoding (i.e., encoding front-to-back and back-to-front information). A random forest classifier (RF) was built with 100 trees and used a voting strategy to integrate these results to predict the interaction.	BILSTM extracted the sequence and position of the biological information in the protein sequence.	
Heterogeneous Network [2021] [296]	PPI prediction was performed using a computational sequence and network representation learning-based model. Local features were extracted from the protein sequence using the $k$ -mer method ( $k = 3$ ), while global features were extracted from the heterogeneous network. The latter captured network structure and obtained potential linked information. This method integrated local features with global features to represent protein nodes.	The protein node contained protein attribute and network structure information by integrating local and global features.	Model accuracy is relatively low compared to other deep learning methods such as DPPI [298].

Framework	Description	Advantage	Disadvantage
OR-RCNN [2021] [301]	<p>This method was called ordinal regression and recurrent convolutional neural network (OR-RCNN), which predicted PPIs based on their confidence score. The architecture comprised two recurrent convolutional neural networks (RCNNs) encoders, which shared the same parameters, to extract robust local features and sequential information from protein pairs.</p> <p>Then, one novel embedding vector was obtained by element-wise multiplication of the two embedding vectors from RCNNs.</p> <p>The second part of the architecture performed an ordinal regression model via multiple sub-classifiers that use the ordinal information behind the confidence score.</p> <p>Finally, the confidence score determined the existence of PPI with a threshold.</p>	<p>This method offered better accuracy compared to some existing models, such as autocovariance [288] and composition transition distribution (CTD) descriptor [308] for feature description, and random forest (RF) [309], extreme gradient boosting (XGBoost) [310], and support vector machine (SVM) [311] for the prediction.</p>	<p>The RCNN was built using bidirectional gated recurrent units (bidirectional-GRU). However, GRUs suffer from slow convergence and low learning efficiency [305].</p>
DeepTrio [2022] [302]	<p>The DeepTrio method used a deep-learning framework based on a mask multiscale CNN architecture that performed binary PPI prediction by capturing multiscale contextual information of protein sequences using multiple parallel filters.</p> <p>This method used a single-protein class, allowing it to distinguish relative and intrinsic properties.</p> <p>This method was also made available as an online tool to address cross-platform usage and dependency-related issues.</p>	<p>DeepTrio is available both online and offline.</p>	<p>DeepTrio yields lower PPI prediction accuracy on the <i>S.cerevisiae</i> core dataset from PIPR based on 5-fold cross-validation compared to PIPR [299].</p> <p>DeepTrio achieves lower PPI prediction accuracy on the <i>S.cerevisiae</i> core dataset from DeepFE-PPI based on 5-fold cross-validation compared to DeepFE-PPI [312].</p>



Contents lists available at ScienceDirect

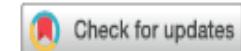
## Journal of Pharmaceutical Sciences

journal homepage: [www.jpharmsci.org](http://www.jpharmsci.org)



### Review

## Prediction Machines: Applied Machine Learning for Therapeutic Protein Design and Development



Tim J. Kamerzell <sup>a, b, \*</sup>, C. Russell Middaugh <sup>a</sup>

<sup>a</sup> Department of Pharmaceutical Chemistry, The University of Kansas, Lawrence, KS, USA

<sup>b</sup> Division of Internal Medicine, HCA MidWest Health, Overland Park, KS, USA

Characteristics of Machine Learning Models Used in Pharmaceutical Protein Development.

	Machine Learning Algorithms and Tasks	Types of Data Input and Features	Perf	Ref
Asparagine Aspartic Acid Degradation	RF Classification and Regression	<ul style="list-style-type: none"> <li>Dataset of n = 776 mAb peptides. Training set includes 64 IgG1s and 3 IgG4s. Validation set contains 10 IgG1s and 2 IgG4s</li> <li>Binary classifier for deamidation lability and regression analysis to predict deamidation rate</li> <li>Features include 12 descriptors for primary sequence (2), backbone orientation (3), side chain orientation (2), solvent accessibility (2) and hydrogen bonding (3)</li> </ul>	Acc 100% R <sup>2</sup> 0.93	[52]
	SVM, RF,NB,KNN,ANN,PLS Classification	<ul style="list-style-type: none"> <li>Training set consists of 194 Asn residues and 25 proteins. Test set consists of 81 Asn residues and 3 protein structures</li> <li>Binary classifier for deamidation lability</li> <li>Descriptors include one experimental measure (deamidation half-life) and several structural features including 3D structure based properties, backbone orientation, side chain orientation, local secondary structure, B-factors, solvent accessibility and reaction coordinate</li> <li>Feature selection with recursive feature elimination</li> <li>10 fold cross validation</li> <li>Data sets unbalanced toward the negative case. Therefore used recall and specificity for further characterization</li> <li>Dataset of 37 mAbs including 55 Asn hotspots and 940 Asn non hotspots, 40 Asp hotspots and 1425 Asp non hotspots.</li> <li>Binary classifier for hotspot</li> <li>Features include 20 descriptors, secondary structure, transition state accessibility, nucleophilicity, pKa, solvent accessibility</li> <li>Dataset unbalanced toward non hotspots therefore a standard weighting scheme using inverse of class frequency employed</li> <li>Monte Carlo cross validation used</li> <li>Dataset of 5 IgG1 and 5 IgG4 mAbs</li> <li>Classification 4 groups with increasing deamidation propensity</li> <li>Features include backbone flexibility, C-N distance and solvent accessibility</li> <li>No clear validation method or performance metrics provided</li> </ul>	Acc 95% AUC 0.96 Recall 0.80 Spec 0.96	[53]
	SVMD,ANN,DT, RF Classification	<ul style="list-style-type: none"> <li>Dataset of 37 mAbs including 55 Asn hotspots and 940 Asn non hotspots, 40 Asp hotspots and 1425 Asp non hotspots.</li> <li>Binary classifier for hotspot</li> <li>Features include 20 descriptors, secondary structure, transition state accessibility, nucleophilicity, pKa, solvent accessibility</li> <li>Dataset unbalanced toward non hotspots therefore a standard weighting scheme using inverse of class frequency employed</li> <li>Monte Carlo cross validation used</li> <li>Dataset of 5 IgG1 and 5 IgG4 mAbs</li> <li>Classification 4 groups with increasing deamidation propensity</li> <li>Features include backbone flexibility, C-N distance and solvent accessibility</li> <li>No clear validation method or performance metrics provided</li> </ul>	TPR 0.94	[52]
	DT Classification	<ul style="list-style-type: none"> <li>Dataset of 5 IgG1 and 5 IgG4 mAbs</li> <li>Classification 4 groups with increasing deamidation propensity</li> <li>Features include backbone flexibility, C-N distance and solvent accessibility</li> <li>No clear validation method or performance metrics provided</li> </ul>		[53]
Protein Oxidation	SVM,RF,NN Classification	<ul style="list-style-type: none"> <li>Dataset of 113 polypeptides containing 975 Met residues of which 122 oxidation prone and 853 oxidation resistant</li> <li>Binary classifier for oxidation</li> <li>54 Features include solvent accessibility, primary structure, met-aromatic residue distance, 3D structure</li> <li>Feature extraction and selection with mRMR and GJ</li> <li>10 fold cross validation</li> <li>Stratified random sampling used due to unbalanced data</li> </ul>	Acc 84% Spec 0.87 F-meas 0.48	[64]
	SVM, RF, NN, IDA Classification	<ul style="list-style-type: none"> <li>16-46 proteins with 2616 methionine sulfoxides</li> <li>Binary for oxidized</li> <li>54 Features including 40 primary structure, 14 tertiary structure, 2 solvent accessibility, 2 entropy and 1 frequency feature</li> <li>Feature extraction and selection with AI and GJ</li> <li>10 fold cross validation</li> <li>Stratified random sampling used due to unbalanced data</li> </ul>	Acc 82% AUC 0.85 Sens 77 Spec 82	[65]
	RF, MLR Classification and Regression	<ul style="list-style-type: none"> <li>172 Met residues across 122 mAbs</li> <li>Binary for oxidized and RF regression model trained to predict relative change in % oxidized species</li> <li>18 descriptors initially evaluated. Final model included 2 structural and 2 dynamic features including the number of overlaps between atoms of Met residue with neighbors, solvent accessibility, and mean square fluctuation of Met C alpha atom</li> <li>Feature extraction and selection with GJ</li> <li>Hold out validation</li> </ul>	R 0.94 Acc 90% AUC 0.96 Spec 0.98 Sens 0.83	[67]
	NN Classification	<ul style="list-style-type: none"> <li>166 peptide fragments of which 32 are positive and 134 negative</li> <li>Binary for oxidized</li> <li>Initially 735 features which was decreased to the top 16 features including secondary structure, solvent accessibility, disorder, amino acid properties and side chain count of carbon atom deviation from mean</li> <li>Feature selection with mRMR</li> <li>Jackknife cross validation</li> </ul>	Acc 92% Sens 0.84 Spec 0.94	[66]
Protein Solution Viscosity	Log linear, PCR Regression	<ul style="list-style-type: none"> <li>14 mAbs at a single high concentration</li> <li>Features include Fv net charge, Fv charge symmetry and Fv hydrophobicity</li> <li>Leave one out cross validation</li> </ul>	R 0.8	[55]
	Log linear, stepwise linear Regression	<ul style="list-style-type: none"> <li>viscosity from concentration dependence of 16 mAbs</li> <li>Both log linear and stepwise linear regression used</li> <li>Features include sequence structural attributes including charge, pI, zeta potential, dipole moment, solvent accessibility and nonpolar surface</li> <li>Slope of log linear equation (B) and stepwise linear approach used to identify correlations and predict B</li> <li>leave one out cross validation</li> </ul>	varies	[54]

T.J. Kumeroff, C.R. McLaughlin / Journal of Pharmaceutical Sciences 110 (2021) 665–681

\* Source: xxxx, 2022.0x.xx

Machine Learning Algorithms and Tasks	Types of Data Input and Features	Perf	Ref
Log linear Regression	<ul style="list-style-type: none"> <li>Viscosity from concentration dependence of 11 mAbs</li> <li>Slope of log linear equation and stepwise linear regression used</li> <li>Features include net charge, zeta potential, isoelectric point, hydrophilic and hydrophobic surface area and Kyte-Doolittle hydrophobicity moment</li> <li>leave one out cross validation</li> </ul>	R <sup>2</sup> 0.88 -0.96	51
ANN, SVM, ensemble, DT, NB, KNN, LDA, PCA Classification and Regression	The authors in this review used the viscosity dependent data and features provided in references <sup>51,53</sup> for both classification and regression of protein solution viscosity. for a detailed discussion of the models and parameters used see manuscript text and figure legends	varies	TJK CRM
Sub-visible Particles KLD, MDS, clustering Classification	<ul style="list-style-type: none"> <li>Classification of particle type from flow microscopy images after freeze thaw stress, shaking, pH stress, thermal stress or formulation</li> <li>3 samples per stress and two replicates per sample measured by flow microscopy</li> <li>number of particles detected by flow microscopy in all the samples ranged from 198 to 45,431 with a median of 21,062</li> <li>Divergence calculated using particle descriptors including size and aspect ratio</li> <li>Matrix of divergences projected using MDS</li> <li>100,000 particles per stress source which includes friability, freeze-thawing, agitation and heating</li> <li>Feature extraction using PCA on 37 texture metrics from particles used to train classifiers</li> <li>images also used as direct inputs in training CNN classifiers</li> <li>n = 600 randomly selected individual particles with a sampling frequency, f= 20</li> <li>10 fold cross validation</li> <li>2 x 10<sup>4</sup> total training samples</li> </ul>	Acc 100%	52,54
CNN, SVM, KNN, DT, ensemble Classification	<ul style="list-style-type: none"> <li>CNN to identify aggregation inducing stress which includes freeze thaw stress, agitation stress, pump recirculation and protein and silicone oil mixtures</li> <li>CNN was three layers and a total of 28,640 trainable parameters</li> <li>Protein particle concentrations &gt;50,000/ml.</li> <li>Random Forest used to classify and count silicone oil and non-silicone oil particles in formulations</li> <li>Features include the aspect ratio, circularity, intensity mean, intensity standard deviation, intensity minimum and intensity maximum</li> <li>Random forest filters compared to Logistic S-factor filters</li> <li>k-fold cross validation</li> </ul>	Acc 95–99%	53
RF, logistic regression Classification	<ul style="list-style-type: none"> <li>Random Forest used to classify acid stable and neutral proteins</li> <li>393 proteins and 889 sequence based features</li> <li>GI used to rank feature importance. Top 10–25 features used</li> <li>five-fold cross validation</li> </ul>	varies	55
Biophysical Stability RF Classification	<ul style="list-style-type: none"> <li>Random Forest used to classify four mAbs and predict concentration</li> <li>Crowdsourced machine learning analyses</li> <li>Collected ~350 spectra for each of four mAbs</li> <li>10 fold cross-validation</li> </ul>	Acc 75% AUC 0.91	106
Bagging, LDA, PCA, PLS Classification and Regression	<ul style="list-style-type: none"> <li>Raman spectroscopy used to classify four mAbs and predict concentration</li> <li>2066 features from 4 analytical methods</li> <li>MI for feature ranking and PCA for visualization</li> <li>Parameter optimization with grid search</li> <li>Leave one and leave 8 out cross-validation</li> </ul>	varies	76
KNN, SVM, LDA, QDA, NB, DT, RF, AdaBoost, MI Classification	<ul style="list-style-type: none"> <li>Bio-similarity assessment of 4 mAbs in 2 different formulations resulting in 8 samples in triplicate for total of 24 samples</li> <li>2066 features from 4 analytical methods</li> <li>MI for feature ranking and PCA for visualization</li> <li>Parameter optimization with grid search</li> <li>Leave one and leave 8 out cross-validation</li> </ul>	Acc 100%	76
ANN, PLS Regression and Classification	<ul style="list-style-type: none"> <li>6 mAbs and 24 conditions per protein resulting in 144 samples</li> <li>Only number of each amino acid species of the protein used as input parameter</li> <li>light scattering and fluorescence used to measure melting temperature (Tm), aggregation (Tagg) and kD</li> <li>ANN and PLS for prediction of Tm, Tagg</li> </ul>	R <sup>2</sup> 0.98 0.94	54
KNN, SVM, DT, AdaBoost, RF, NB, LDA, PCA, MI Classification	<ul style="list-style-type: none"> <li>35 samples</li> <li>Multiple biophysical techniques including UV–Vis absorption spectroscopy, FTIR, CD and HPLC used to characterize Crofelemer</li> <li>Approximately 7000 features per sample</li> <li>PCA for feature scaling and MI to identify important features</li> <li>Monte Carlo cross-validation</li> </ul>	Acc 99%	52

\* Source: xxxx, 2022.0x.xx



OXFORD

Briefings in Bioinformatics, 2022, 23(4), 1–20

<https://doi.org/10.1093/bib/bbac267>

Advance access publication date: 14 July 2022

Review

# Machine-designed biotherapeutics: opportunities, feasibility and advantages of deep learning in computational antibody discovery

Wiktoria Wilman, Sonia Wróbel, Weronika Bielska, Piotr Deszynski, Paweł Dudzic, Igor Jaszczyszyn, Jędrzej Kaniewski,

Jakub Młokosiewicz, Anahita Rouyan, Tadeusz Satława, Sandeep Kumar , Victor Greiff  and Konrad Krawczyk 

Corresponding author: Dr Konrad Krawczyk, NaturalAntibody. E-mail: konrad@naturalantibody.com, konrad@naturalantibody.com

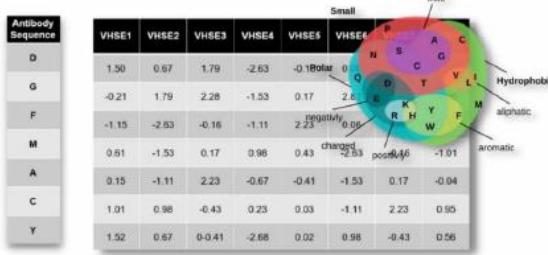
### A. One-hot encoding

Antibody Sequence	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
D	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
G	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
F	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
M	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	
A	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
C	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
Y	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	

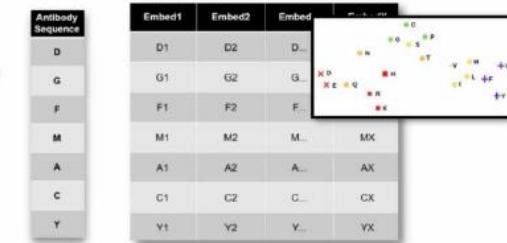
### B. Substitution matrix-based

Antibody Sequence	1	2	3	4	5	6	D	G	F	M	A	C	Y
D	0.01	0.04	0.00	0.03	0.05	0.01	-7	-3	-3	-3	-1	-2	-2
G	-0.3	0.9	-0.1	-0.3	0.2	0.0	-3	9	-1	-3	-6	1	-1
F	-3	-1	9	2	2	0	2	-1	10	-7	-5	0	-1
M	-3	-3	2	10	-7	-7	-3	-3	-7	-7	-1	2	-1
A	-1	-6	-5	-7	13	-7	-7	-7	-7	-7	-6	-6	-6
C	-2	1	0	-1	-1	9	-1	-1	-1	-1	9	4	4
Y	-2	-1	-1	2	-6	4	-6	-6	-6	-6	4	8	8

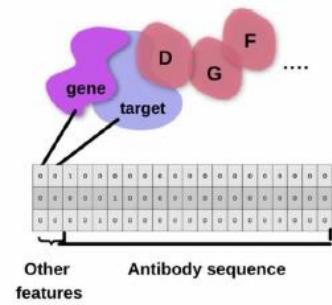
### C. Amino acids properties



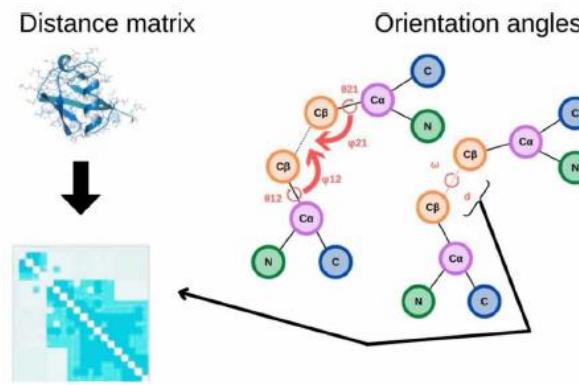
### D. Learned amino acids properties



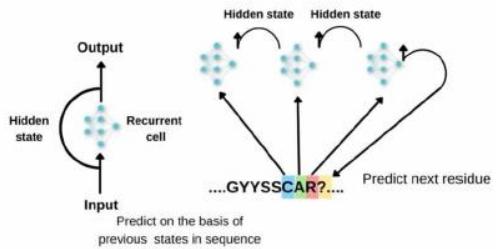
### E. Encoding of supplementary attributes



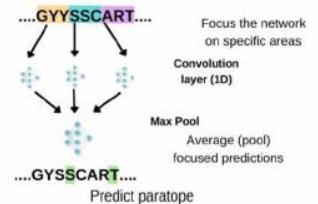
### F. Encoding of structural features



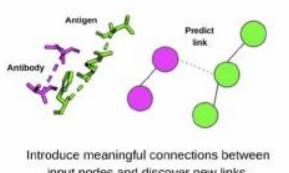
### A. Recurrent Neural Networks



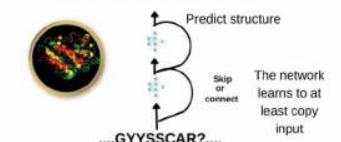
### B. Convolutional Neural Networks



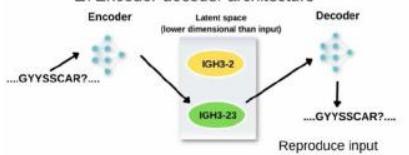
### C. Graph Neural Networks



### D. Residual Neural Network



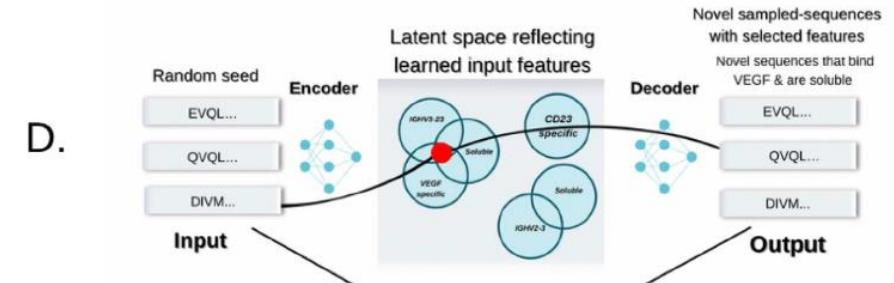
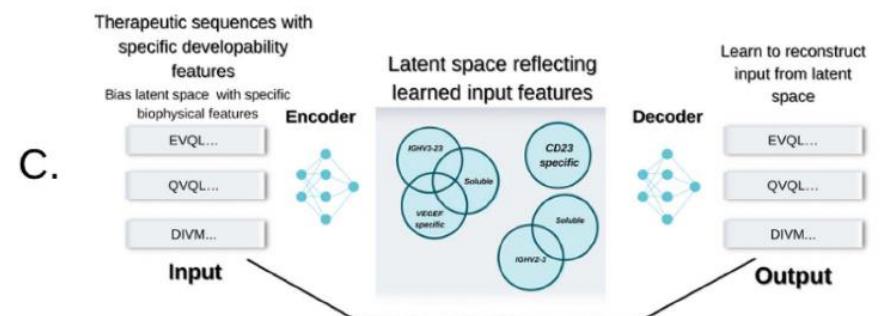
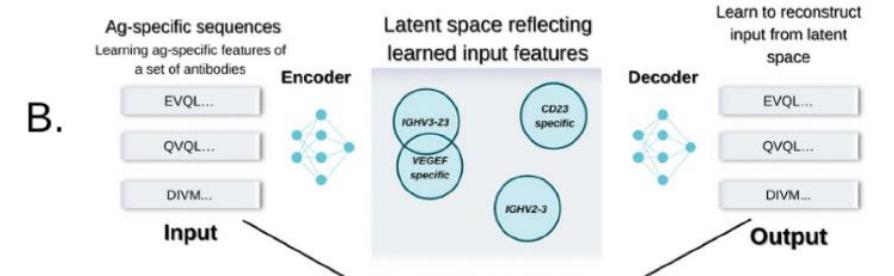
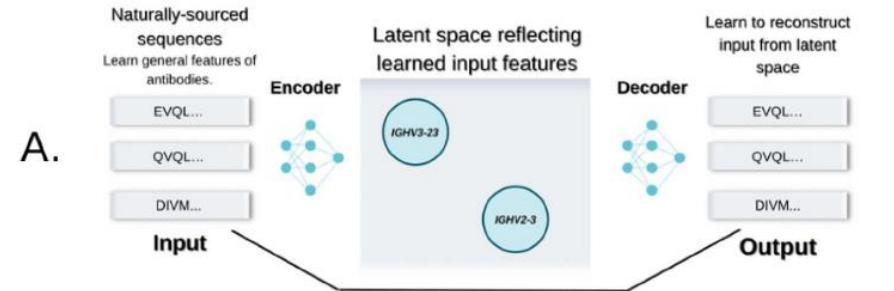
### E. Encoder-decoder architecture



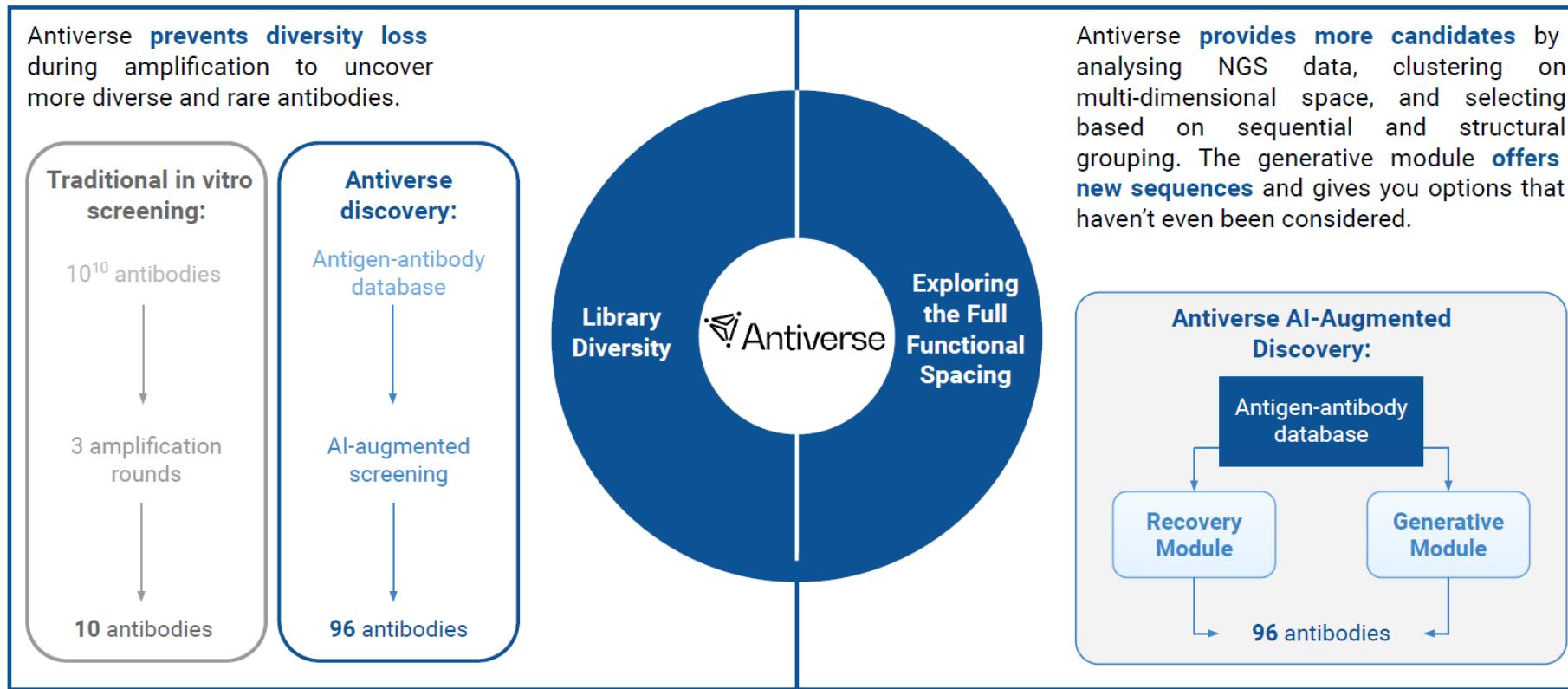
Category	Method	Problem solved	Training input	Architecture	Training parameters	Libraries	Availability	Paper
Structure prediction	DeepH3	CDRH3 prediction	1388 structures	Series of 1D and 2D convolutions (3 1D + 25 2D blocks)	30 epochs, batch size 4, 35 h using one NVIDIA Tesla K80 Graphics processing unit (GPU)	PyTorch	<a href="#">link</a>	[27]
	DeepAb	V region structure prediction	118 386 sequences and 1692 structures	A 1D ResNet (1D convolution followed by three 1D ResNet blocks) and the bi-LSTM encoder	60 epochs, batch size 128, NVIDIA K80 GPU requiring 60 h	PyTorch	<a href="#">link</a>	[19]
	AbLooper	CDR Prediction	3438 structures	Five E(n)-equivariant graph neural networks (EGNNs), each one with four layers	NVIDIA Tesla V100 GPU, predict the CDRs for one hundred structures in under five seconds	PyTorch	<a href="#">link</a>	[40]
	NanoNet	Heavy chain prediction	~2000 structures	Two 1D ResNets with input tensor of 140 × 22	batch size of 16 and ~130 epochs, 10 min on a GeForce RTX 2080 Ti	Keras/TensorFlow	<a href="#">link</a>	[18]*
Humanization/Deimmunization	Nativeness LSTM	Learn distribution of amino acids at positions	400 000 sequences	Bidirectional LSTM with dimensionality 64	10 epochs	PyTorch	<a href="#">link</a>	[31]
	Sapiens	Antibody humanization	20 million heavy chains and 19 million light chains	RoBERTA transformer, 4 layers, 8 attention heads, 568 857 parameters	700 epochs for heavy chains, 300 epochs for light chains	PyTorch-/Fairseq [41]	<a href="#">link</a>	[24]
	hu-Mab	Discriminate between human/mouse sequences	65 million sequences with 13 million non-human ones	Random Forest	n/a	scikit-learn	<a href="#">link</a>	[42]
Binding models	Parapred	Paratope residues prediction	1662 sequences (277 antibody-antigen complexes × 6 Complementarity determining regions each) and tested on the same dataset using 10-fold cross-validation technique	Convolutional and recurrent neural networks	16 epochs, 32 batch size	Keras	<a href="#">link</a>	[43]
	Epitope3d	Conformational epitopes prediction	1351 antibody-antigen structures (covering 40 842 epitope residues) and 180 unbound antigen structures; tested on 20 unbound antigen structures; 45 unbound antigen structures used for external blind test	Supervised learning algorithms: Multi-layer Perceptron, Support Vector Machines, K-Nearest Neighbor, AdaBoost, Gaussian processes (GP), Random Forest, Gradient Boost, XGBoost, Extra Trees	n/a	scikit-learn Python	<a href="#">link</a>	[44]
	mmCSM-AB	Prediction of the consequences of multiple point mutations on antibody-antigen binding affinity	1640 mutations with associated changes in binding affinity (905 single missense mutations and 735 modeled reverse mutations); tested on 242 multiple missense mutations with associated changes in binding affinity	Supervised learning algorithms for example: Random Forest, Extra Trees, Gradient Boost, XGBoost, SVM and Gaussian Process	n/a	scikit-learn Python	<a href="#">link</a>	[45]

Category	Method	Problem solved	Training input	Architecture	Training parameters	Libraries	Availability	Paper
Sequence-based methods	Phage display LSTM	Generate novel kynurenine binding sequences from LSTM	959 sequences	LSTM, two layers with 64 units.	269 epochs	Keras/Tensorflow	n/a	[46]
	Phage display CNN	Predict phage enrichment and generate novel CDRH3	96 847 sequences (largest dataset on github)	Ensemble of CNNs, largest with two convolutional layers and 18 706 parameters	20 epochs	Keras	<a href="#">link</a>	[47]
	Image-based prediction	Distinguish between binding antibodies and lineages	24 953 models with calculated fingerprints from 308 EBOV and 54 HIV antibodies.	ResNet-50 [48]	Pre-trained model	Keras/Tensorflow	<a href="#">link</a>	[33]
	Paratope and Epitope Prediction with graph Convolution Attention Network (PECAN)	Epitope and paratope prediction	162 structures for epitope prediction and 460 for paratope prediction	Graph Convolutional Attention Network	Up to 250 epochs, batch size of 32 (multiple parameters tested)	Tensorflow	<a href="#">link</a>	[20]
	DLAB	Sorting of protein docking poses	759 Antibody-antigen complexes	Convolutional Neural Network	n/a	PyTorch	<a href="#">link</a>	[32]
	immune2vec	Embed CDRH3 into 100 dimensions using skip-gram	15,63 million sequences	Two dense layers	n/a	Gensim	<a href="#">link</a>	[49]
	ProtVec	CDRH3 sequences to predict COVID-19 status	COVID-19 data from OAS	Based on ProtVec from Harvard DataVerse [50] and SVM	Reused previous model.	Reused previous model	<a href="#">link</a>	[51]*
	AntiBerty	Masked language modeling, paratope prediction	558 million sequences	BERT transformer encoder model, 8 layers, 26 M trainable A100 GPUs parameters.	8 epochs, 10 days on four NVIDIA	PyTorch	n/a	[39]*
	AntiBerta	Masked language modeling, paratope prediction	57 million sequences	Antibody-specific Bi-directional Encoder Representation from Transformers, 86 m parameters	12-layer transformer model that is pre-trained on 57 M human BCR sequences, 3 epochs, batch size of 96 across 8 NVIDIA V100 GPUs	PyTorch	n/a	[25]
	AbLang	Masked language modeling, reconstruct erroneous sequences	14 million heavy chains, 200 000 light chains training. Evaluation sets of 100 k, 50 k for heavy lights respectively.	Based on RoBERTA from HuggingFace. 12 layers.	20 epochs for heavy chains, batch 8192, light chains 40 epochs 4096 batch size	PyTorch	<a href="#">link</a>	[52]*
Generative methods/anti-body design	Mouse VAE	Model latent space of CDR triples of antigen challenged mice	243 374 sequences.	VAE with encoder and decoder each having two dense layers (256 512 units each)	200 Epochs on a single GPU from the ETH cluster.	Tensorflow	n/a (available after peer review)	[53]*
	Developability controlled GAN	Learn latent representation of human sequences and bias it towards biophysical properties	400 000 sequences	Generative Adversarial Network, size of 128 (single chain) seven layers consisting of 2D convolution and dense layers.	500 epochs, batch size of 128	Keras/Tensorflow	n/a	[35]*
	Nanobody generation	Autoregression on nanobody sequences to generate novel CDRH3	1.2 million sequences	ResNet with nine blocks with six dilated convolutional layers.	250 000 updates, batch size of 30.	Tensorflow/PyTorch	<a href="#">link</a>	[17]

Category	Method	Problem solved	Training input	Architecture	Training parameters	Libraries	Availability	Paper
	In silico LSTM	In silico proof-of-principle of virtually unconstrained antigen-specific antibody sequence generation	70 000 murine CDR3 sequences	1024 LSTM with embedding layer and dense output layer.	20 epochs, batch size 64	Tensorflow	<a href="#">link</a>	[54]
	Immunoglobulin Language Model (IgLm)	Masked language modeling, generate synthetic libraries of antibodies by solving masked language model	558 million sequences	Transformer decoder architecture based on the GPT-2 model with 512 embeddings, 12 million parameters VAE	batch size of 512 and 2 gradient accumulation steps using DeepSpeed, 3 days when distributed across 4 NVIDIA A100 GPUs	GPT-2 from HuggingFace	n/a	[55]*
	IG-VAE	Immunoglobulins structure generation	10 768 immunoglobulins structures (including 4154 non-sequence-redundant structures)- set covers almost 100% of the antibody structure database (AbDb); Tested on 5000 structures from the latent space of the Ig-VAE	n/a	PyTorch	n/a	[34]*	
	Generative method Benchmarking: (AR) the sequence-based based on portion of autoregressive generative model, geometric vector perceptron (GVP) the precise structure-based graph neural network and (Fold2Seq) fold-based generative model GNN-based generation	Antibody CDR regions design	Sequences from natural llama nanobody repertoire	AR- Autoregressive Causal Dilated Convolutions; GVP-based Encoder-Decoder GNN; Fold2Seq- Encoder-Decoder Transformer	n/a	n/a	n/a	[56]*
	AntBO	CDR3 sequence and 3D structure design	~5000 structures. For CDR-H1, the train/validation/test size is 4050, 359 and 326. For CDR-H2, the train/validation/test size is 3876, 483 and 376. For CDRH3, the train/validation/test size is 3896, 403 and 437.	Message passing network (MPN): Iterative Refinement Graph Neural Network (RefineGNN)	batch size of 16, dropout of 0.2 and learning rate of 0.0005	n/a	n/a	[23]*
		CDR3 region design		Bayesian Optimization and GP	87 cores 12 GB GPU memory	GPyTorch, Botorch	n/a	[57]*



**Antiverse** is a new type of antibody discovery company accelerating drug development. The Antiverse platform exists at the intersection of structural biology, machine learning and medicine to enable breakthroughs to happen more quickly and cost-effectively.

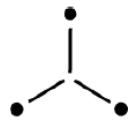




**Antiverse** is recognized as one of the top biotech startups in the UK with our antibody discovery service already in use by big pharma. The main feature of the company is **10x Diversity with AI-Augmented Drug Discovery**.

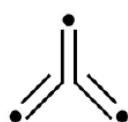
**Existing antibody** discovery methods are well-developed and often effective at discovering binders. But when there is a need to find the best possible candidate, or when finding a suitable candidate is hard with current methods, the options are **limited** and often **costly**.

Antiverse uses **next-generation sequencing (NGS)** to extract more data from existing workloads. The **AI-Augmented Drug Discovery platform** and trained models analyse the statistics gained from thousands of experiments. These outputs are compared against known data in order to select best candidates.



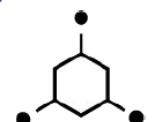
### Target Selection

Antiverse provides targeted options in order to focus on testing safely once there are too many antibody-antigen binding options.



### Binder Recovery

Antiverse can help to find sufficient potential binders that can be missed by conventional methods.

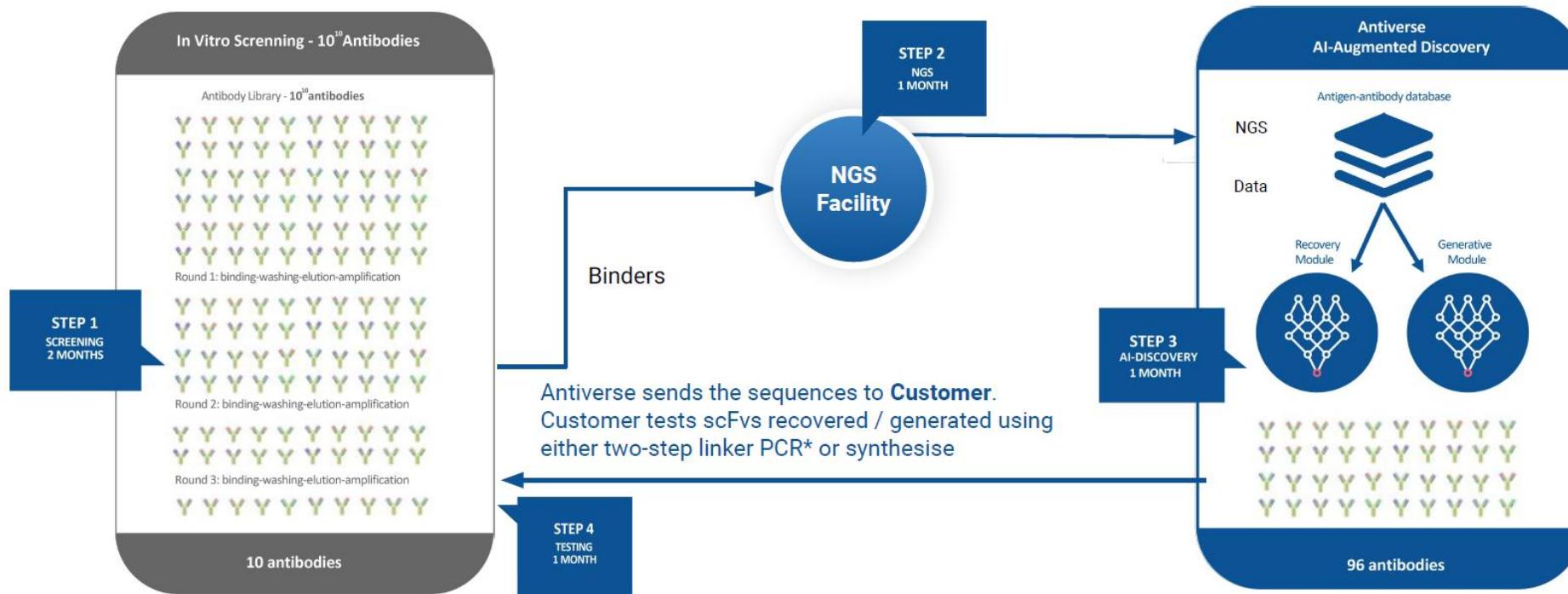


### Binder Customisation

Antiverse can generate new binder variants that will be sufficient for clients purposes.



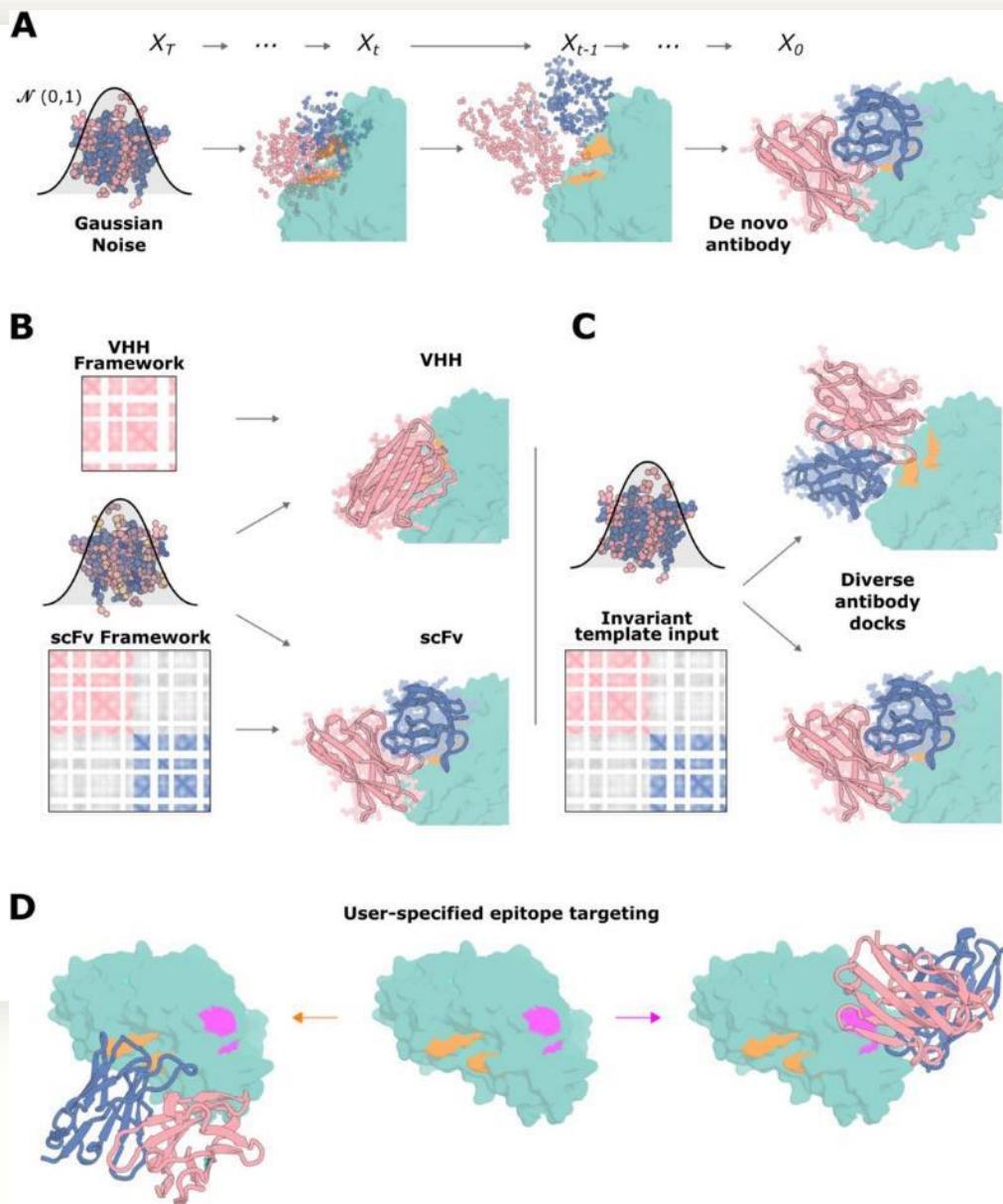
The **Antiverse AI-ADD** system found each and every cluster identified by other methods, plus more. These additional clusters contained rare and unique sequences.



bioRxiv preprint doi: <https://doi.org/10.1101/2024.03.14.585103>; this version posted March 18, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

# Atomically accurate de novo design of single-domain antibodies

Nathaniel R. Bennett<sup>‡1,2,3</sup>, Joseph L. Watson\*<sup>‡1,2</sup>, Robert J. Ragotte<sup>‡1,2</sup>, Andrew J. Borst<sup>‡1,2</sup>, Déjenaé L. See<sup>#1,2,4</sup>, Connor Weidle<sup>‡1,2</sup>, Riti Biswas<sup>1,2,3</sup>, Ellen L. Shrock<sup>1,2</sup>, Philip J. Y. Leung<sup>1,2,3</sup>, Buwei Huang<sup>1,2,4</sup>, Inna Goreshnik<sup>1,2,5</sup>, Russell Ault<sup>6,7</sup>, Kenneth D. Carr<sup>2</sup>, Benedikt Singer<sup>1,2</sup>, Cameron Criswell<sup>1,2</sup>, Dionne Vafeados<sup>2</sup>, Mariana Garcia Sanchez<sup>2</sup>, Ho Min Kim<sup>8,9</sup>, Susana Vázquez Torres<sup>1,2,10</sup>, Sidney Chan<sup>2</sup>, David Baker\*<sup>1,2,5</sup>



NEWS | 19 March 2024

# 'A landmark moment': scientists use AI to design antibodies from scratch

Modified protein-design tool could make it easier to tackle challenging drug targets – but AI antibodies are still a long way from reaching the clinic.

By [Ewen Callaway](#)

Researchers have used generative artificial intelligence (AI) to help them make completely new antibodies for the first time.



BIOTECH

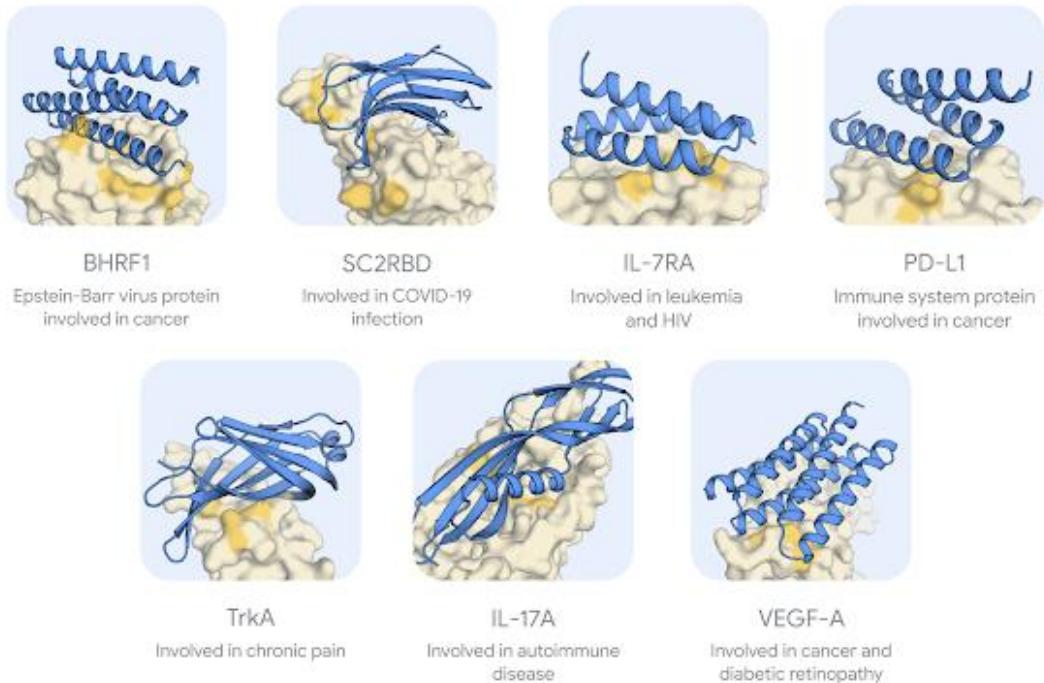
## New AI drug discovery powerhouse Xaira rises with \$1B in funding

By [Annalee Armstrong](#) • Apr 24, 2024 6:00am[ARCH Venture Partners](#) [Artificial Intelligence](#) [drug discovery](#) [Marc Tessier-Lavigne](#)

# *De novo* design of high-affinity protein binders with AlphaProteo

Vinicius Zambaldi<sup>\*,†</sup>, David La<sup>\*,†</sup>, Alexander E. Chu<sup>\*,†</sup>, Harshnira Patani<sup>\*,†</sup>, Amy E. Danson<sup>\*,†</sup>, Tristan O. C. Kwan<sup>\*,†</sup>, Thomas Frerix<sup>\*,†</sup>, Rosalia G. Schneider<sup>\*,†</sup>, David Saxton<sup>\*,†</sup>, Ashok Thillaisundaram<sup>\*,†</sup>, Zachary Wu<sup>\*,†</sup>, Isabel Moraes<sup>2</sup>, Oskar Lange<sup>2</sup>, Eliseo Papa<sup>1</sup>, Gabriella Stanton<sup>1</sup>, Victor Martin<sup>1</sup>, Sukhdeep Singh<sup>1</sup>, Lai H. Wong<sup>1</sup>, Russ Bates<sup>2</sup>, Simon A. Kohl<sup>2</sup>, Josh Abramson<sup>1</sup>, Andrew W. Senior<sup>1</sup>, Yilmaz Alguel<sup>3</sup>, Mary Y. Wu<sup>4</sup>, Irene M. Aspalter<sup>5</sup>, Katie Bentley<sup>5,6</sup>, David L.V. Bauer<sup>7</sup>, Peter Cherepanov<sup>3</sup>, Demis Hassabis<sup>1</sup>, Pushmeet Kohli<sup>1</sup>, Rob Fergus<sup>1,†</sup> and Jue Wang<sup>1,†</sup>

\*Equal contributions, <sup>†</sup>Equal supervision, <sup>1</sup>Google DeepMind, <sup>2</sup>Work performed while at Google DeepMind, <sup>3</sup>The Chromatin Structure and Mobile DNA Laboratory, The Francis Crick Institute, London, UK, <sup>4</sup>COVID Surveillance Unit, The Francis Crick Institute, London, UK, <sup>5</sup>Cellular Adaptive Behaviour Laboratory, The Francis Crick Institute, London, UK., <sup>6</sup>Department of Informatics, King's College London, London, UK. K.B. performed the work at the Cellular Adaptive Behaviour Laboratory, The Francis Crick Institute, London, UK, <sup>7</sup>RNA Virus Replication Laboratory, The Francis Crick Institute, London, UK



## Advancing Targeted Protein Degradation via Multiomics Profiling and Artificial Intelligence

Miquel Duran-Frigola,\* Marko Cigler, and Georg E. Winter\*



Cite This: *J. Am. Chem. Soc.* 2023, 145, 2711–2732



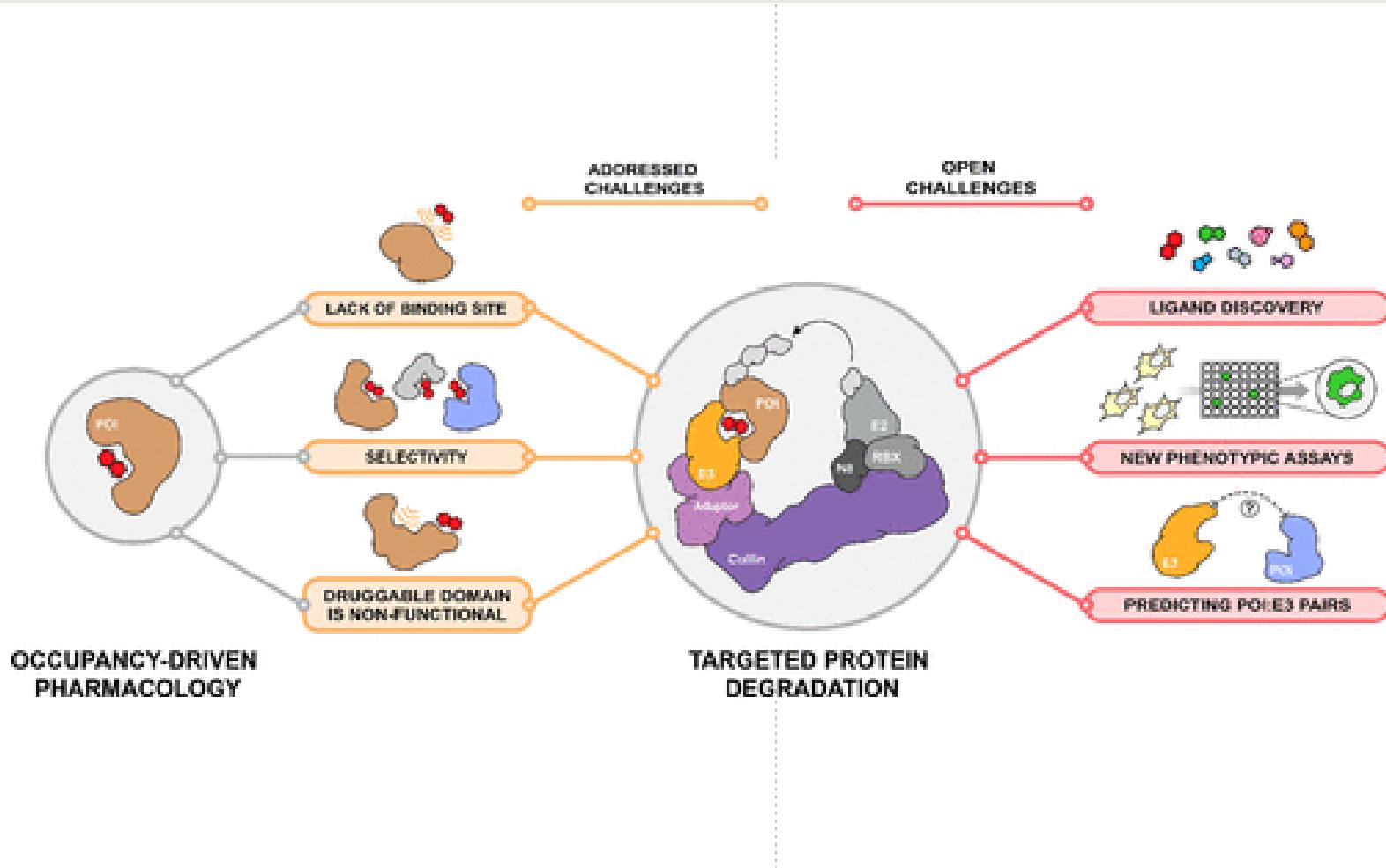
Read Online

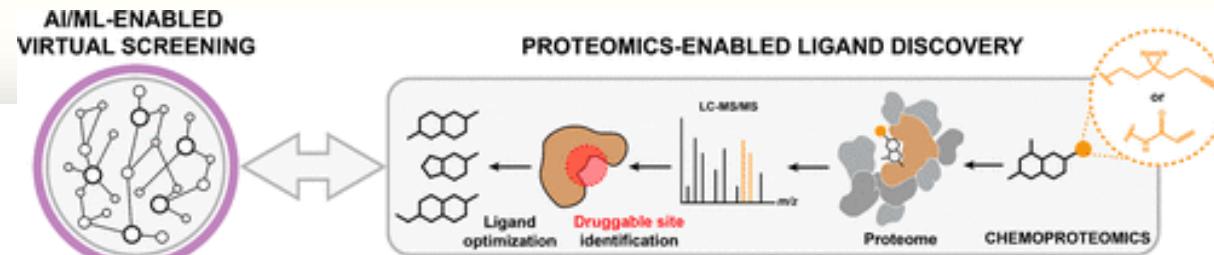
ACCESS |

Metrics & More

Article Recommendations

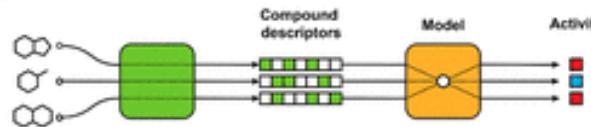
**ABSTRACT:** Only around 20% of the human proteome is considered to be druggable with small-molecule antagonists. This leaves some of the most compelling therapeutic targets outside the reach of ligand discovery. The concept of targeted protein degradation (TPD) promises to overcome some of these limitations. In brief, TPD is dependent on small molecules that induce the proximity between a protein of interest (POI) and an E3 ubiquitin ligase, causing ubiquitination and degradation of the POI. In this perspective, we want to reflect on current challenges in the field, and discuss how advances in multiomics profiling, artificial intelligence, and machine learning (AI/ML) will be vital in overcoming them. The presented roadmap is discussed in the context of small-molecule degraders but is equally applicable for other emerging proximity-inducing modalities.



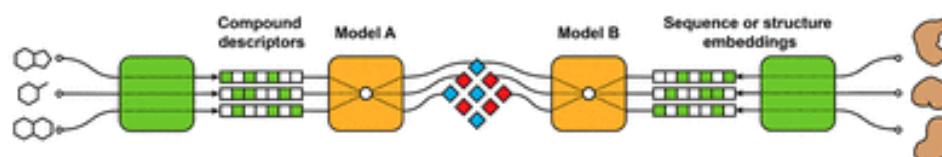


### I METHODS THAT REQUIRE PRE-EXISTING EVIDENCE OF ACTIVITY

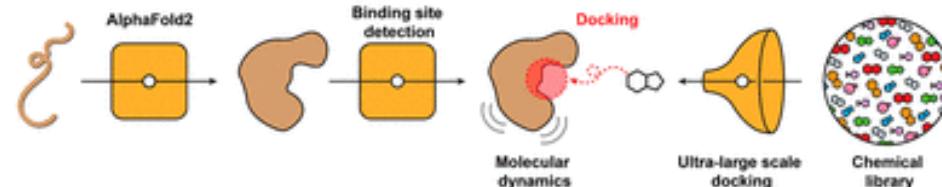
#### QUANTITATIVE SAR



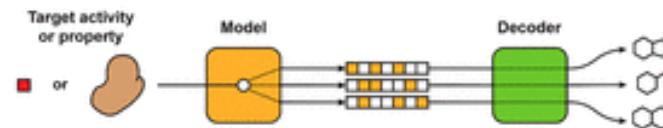
#### PROTEOCHEMOMETRIC

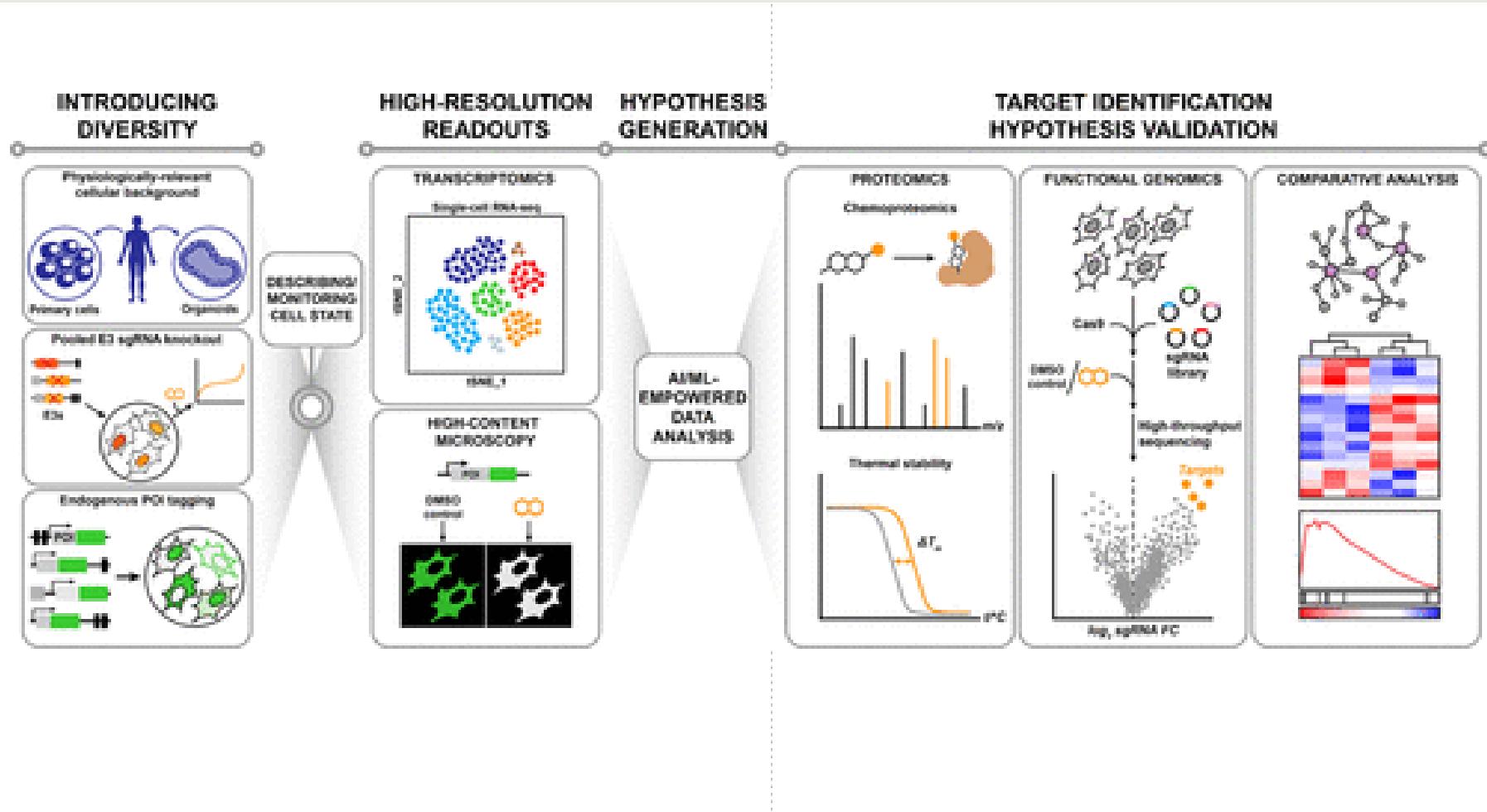


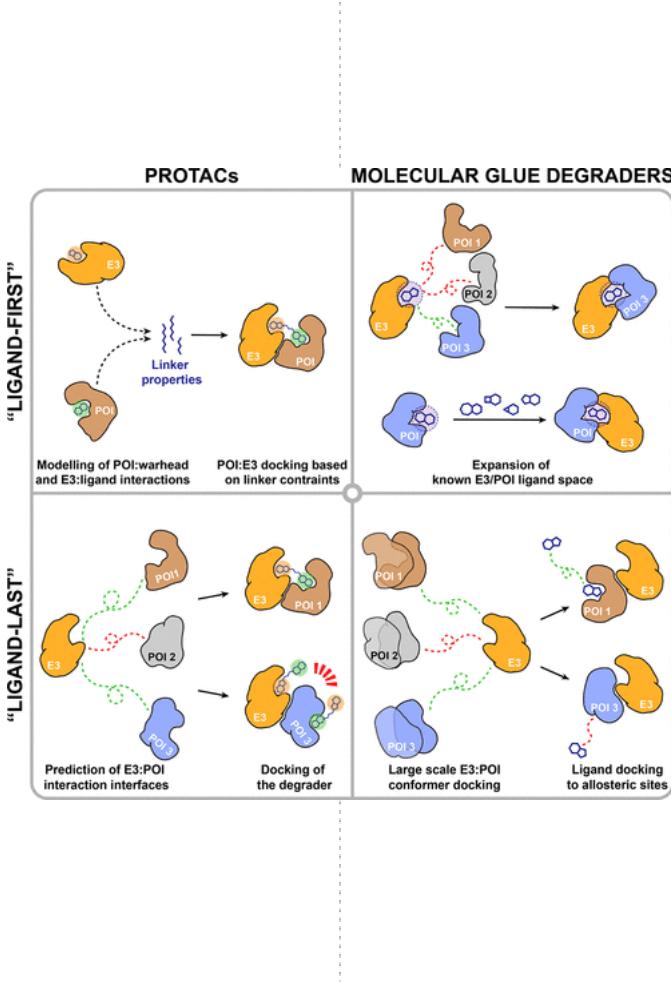
### II PHYSICS-BASED METHODS ENHANCED WITH AI/ML



### III GENERATIVE MODELS









Article

<https://doi.org/10.1038/s41467-022-34807-3>

# DeepPROTACs is a deep learning-based targeted degradation predictor for PROTACs

---

Received: 18 October 2021

---

Accepted: 8 November 2022

---

Published online: 21 November 2022

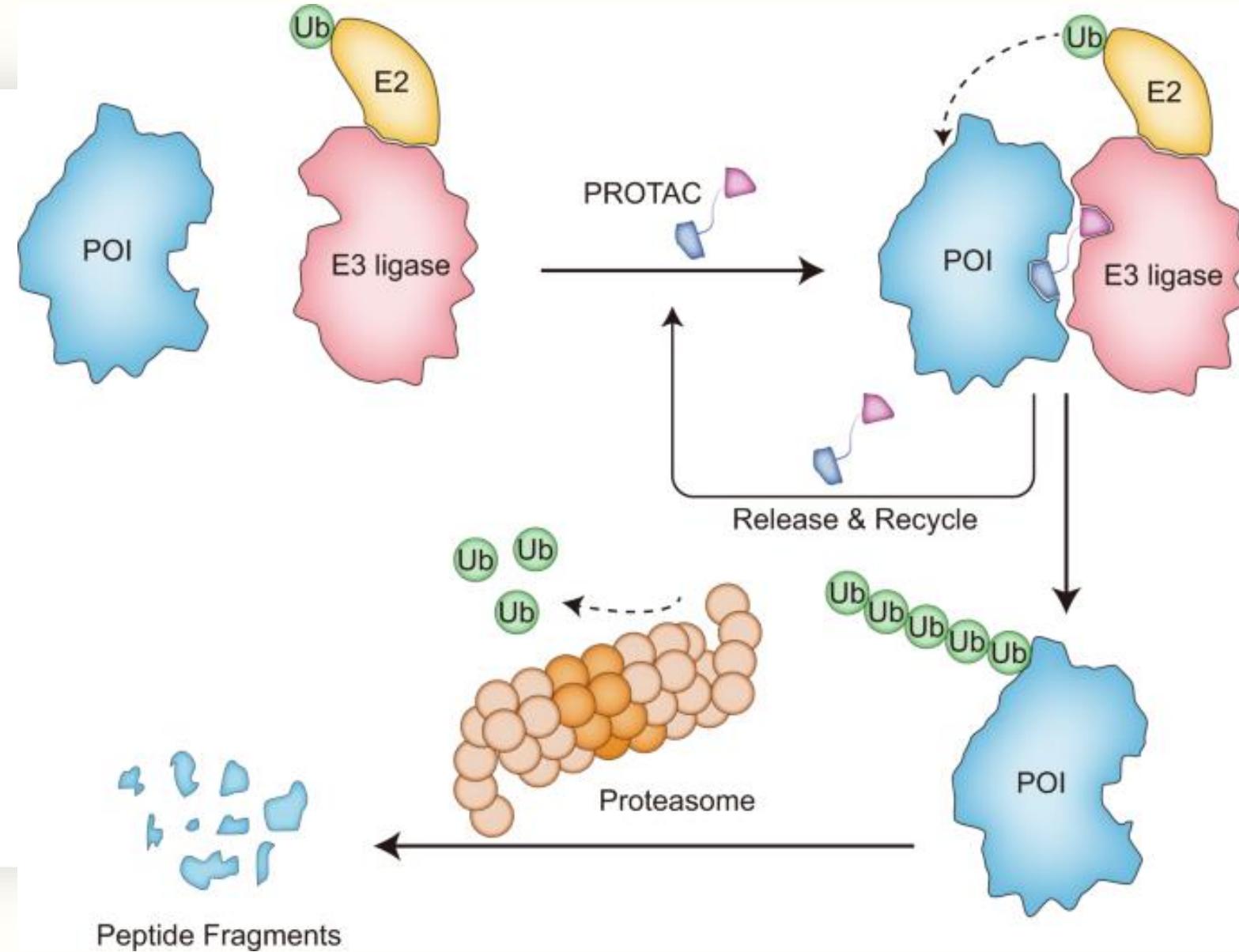
---

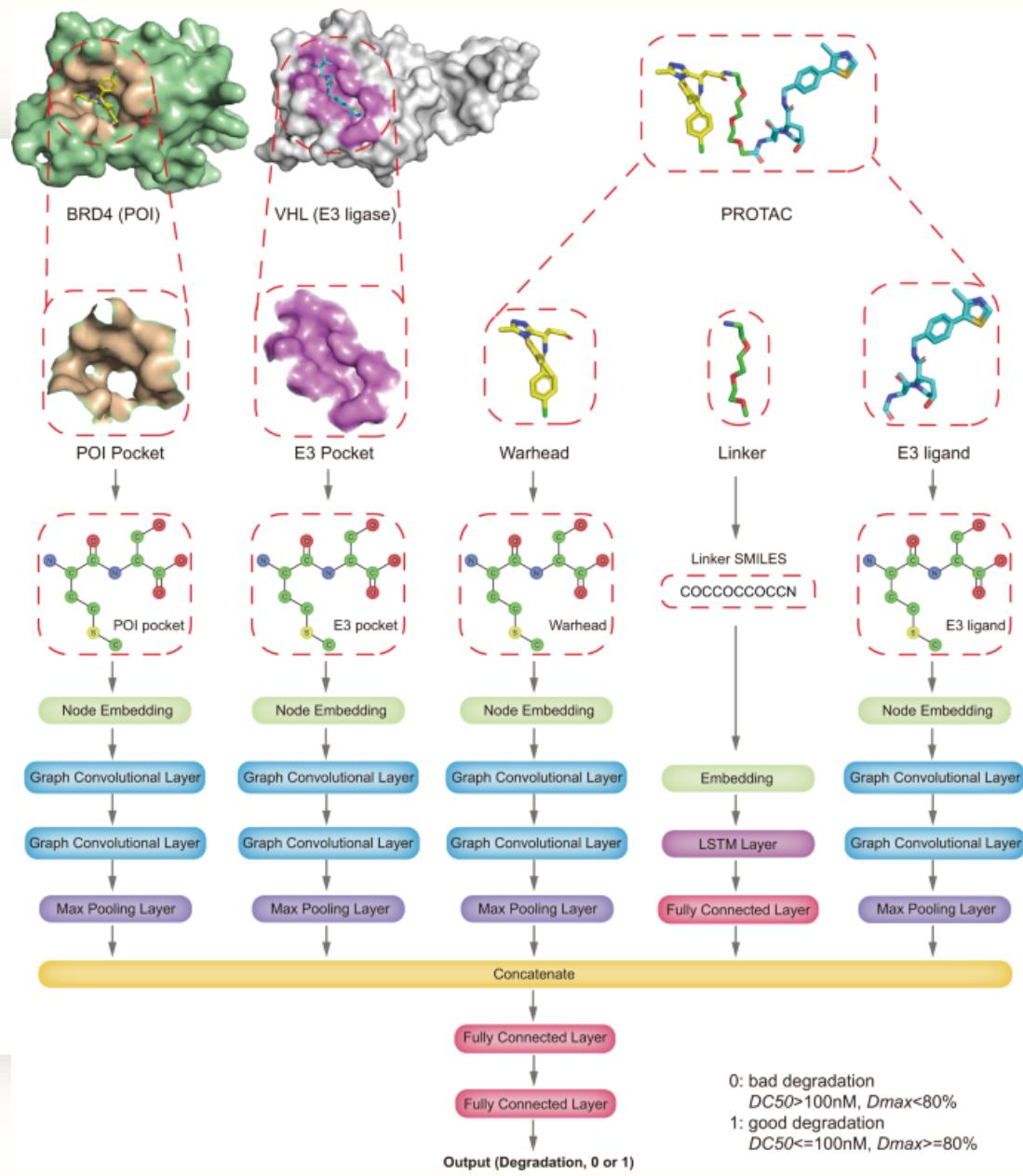
Check for updates

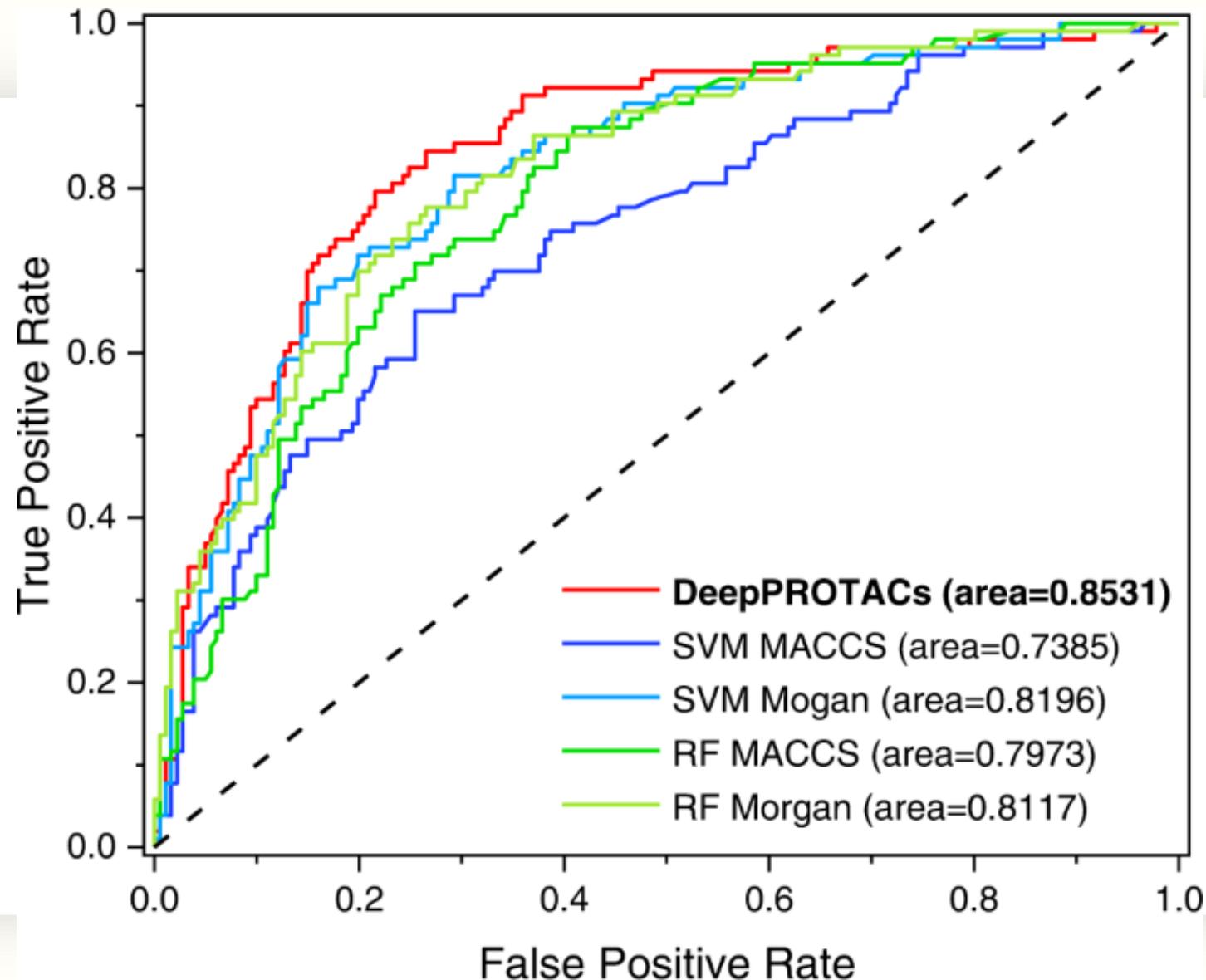
Fenglei Li<sup>1,2,6</sup>, Qiaoyu Hu<sup>1,6</sup>, Xianglei Zhang <sup>1,6</sup>, Renhong Sun<sup>3,6</sup>, Zhuanghua Liu<sup>2,6</sup>, Sanan Wu<sup>1</sup>, Siyuan Tian<sup>1,4</sup>, Xinyue Ma<sup>1,2</sup>, Zhizhuo Dai<sup>4</sup>, Xiaobao Yang <sup>3</sup>✉, Shenghua Gao <sup>2</sup>✉ & Fang Bai <sup>1,2,4,5</sup>✉

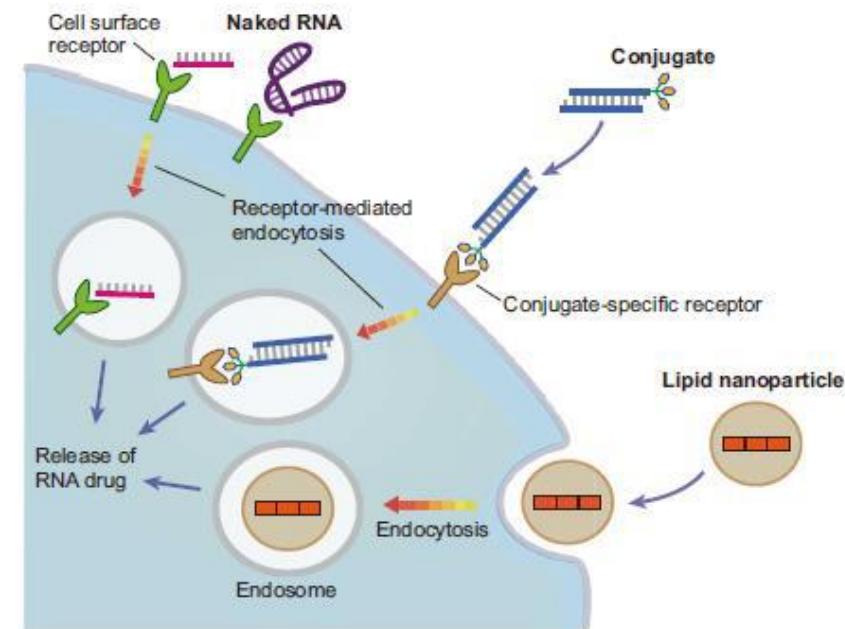
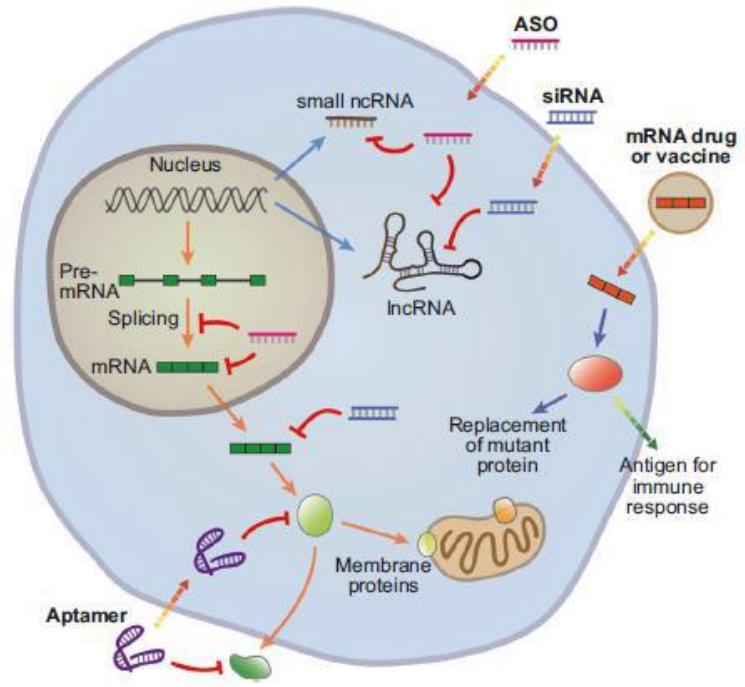
---

The rational design of PROTACs is difficult due to their obscure structure-activity relationship. This study introduces a deep neural network model - DeepPROTACs to help design potent PROTACs molecules. It can predict the degradation capacity of a proposed PROTAC molecule based on structures of given target protein and E3 ligase. The experimental dataset is mainly collected from PROTAC-DB and appropriately labeled according to the  $DC_{50}$  and  $D_{max}$  values. In the model of DeepPROTACs, the ligands as well as the ligand binding









\* Source: xxxx, 2022.0x.xx



# The Limitless Future of RNA Therapeutics

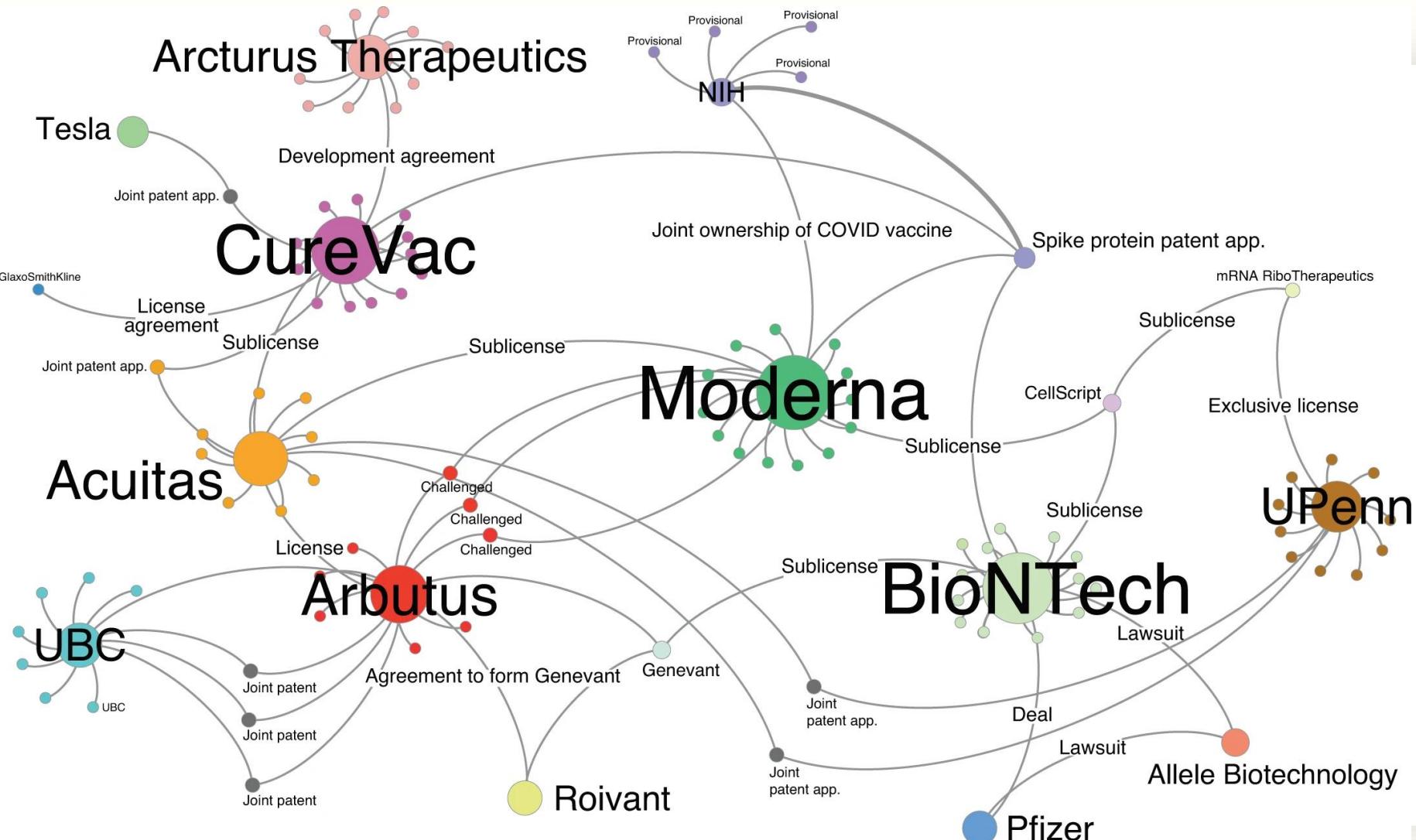
*Tulsi Ram Damase<sup>1</sup>, Roman Sukhovershin<sup>1</sup>, Christian Boada<sup>2</sup>, Francesca Taraballi<sup>3,4</sup>, Roderic I. Pettigrew<sup>2</sup> and John P. Cooke<sup>1\*</sup>*

<sup>1</sup> RNA Therapeutics Program, Department of Cardiovascular Sciences, Houston Methodist Research Institute, Houston, TX, United States, <sup>2</sup> Colleges of Medicine, Engineering, Texas A&M University and Houston Methodist Hospital, Houston, TX, United States, <sup>3</sup> Center for Musculoskeletal Regeneration, Houston Methodist Research Institute, Houston, TX, United States, <sup>4</sup> Department of Orthopedics and Sports Medicine, Houston Methodist Hospital, Houston, TX, United States

**TABLE 1 |** RNA therapeutics approved for clinical use or undergoing clinical trials.

Drug	Type of RNA	Company	Route of delivery	Condition/Disease	Status
Nusinersen (Spinraza)	ASO	Ionis	Intrathecal	Spinal muscular atrophy	FDA approval in 2016
Eteplirsen (Exondys 51)	ASO	Sarepta	Intravenous	Duchenne muscular dystrophy	FDA approval in 2016
Inotersen (Tegsedi)	ASO	Ionis	Subcutaneous	Familial amyloid polyneuropathy	FDA approval in 2018
Volanesorsen (Waylivra)	ASO	Ionis	Subcutaneous	Familial chylomicronemia syndrome	EU approval in 2019
Patisiran (Onpattro)	siRNA	Alnylam	Intravenous	Polyneuropathy	FDA approval in 2018
Givosiran (Givlaari)	siRNA	Alnylam	Subcutaneous	Acute hepatic porphyria	FDA approval in 2019
Cobomarsen (MRG-106)	miRNA	miRagen (Viridian)	Intravenous/subcutaneous	Blood cancers	Phase II
Remlarsen (MRG-201)	miRNA	miRagen (Viridian)	Intradermal	Keloids	Phase II
MRG-110	miRNA	miRagen (Viridian)	Intradermal	Tissue Repair	Phase I
Pegaptanib (Macugen)	Aptamer(RNA)	Bausch + Lomb	Intravitreal	Macular Degeneration	FDA approval in 2014
Emapticap pegol (NOX-E36)	Aptamer(RNA)	NOXXON	Intravenous/Subcutaneous	Diabetic nephropathy, lung and pancreatic cancer	Phase II
Olaptesed pegol (NOX-A12)	Aptamer(RNA)	NOXXON	Intravenous	Brain cancer	Phase I/II
BNT162b2	mRNA	BioNTech and Pfizer	Intramuscular	COVID-19	FDA authorization for emergency use in 2020
mRNA-1273	mRNA	Moderna	Intramuscular	COVID-19	FDA authorization for emergency use in 2020
CVnCoV	mRNA	CureVac	Intramuscular	COVID-19	Phase III
AZD8601	mRNA	Moderna/AstraZeneca	Epicardial	Ischemic heart disease	Phase II
mRNA-1647	mRNA	Moderna	Intramuscular	Cytomegalovirus infection	Phase II
P-BCMA-101	mRNA	Poseida	Intravenous	Multiple myeloma	Phase II
mRNA-4157	mRNA	Moderna	Intramuscular	Cancer	Phase II
mRNA-3704	mRNA	Moderna	Intravenous	Methylmalonic aciduria	Phase I/II
MRT5005	mRNA	Translate Bio	Inhalation	Cystic Fibrosis	Phase I/II
mRNA-2416	mRNA	Moderna	Intratumoral	Solid tumors/lymphoma/advanced ovarian carcinoma	Phase I/II
BNT131 (SAR441000)	mRNA	BioNTech/ Sanofi/Genmab	Intratumoral	Advanced melanoma	Phase I/II
Descartes-08	mRNA	Cartesian	Intravenous	Generalized myasthenia gravis	Phase I/II
BNT122	mRNA	BioNTech	Intravenous	Solid tumors	Phase I/II
mRNA-2752	mRNA	Moderna	Intratumoral	Solid tumors	Phase I
MEDI1191	mRNA	Moderna	Intratumoral	Solid tumors	Phase I
mRNA-1944	mRNA	Moderna	Intravenous	Chikungunya infection	Phase I
CV8102	mRNA	CureVac	Intratumoral	Solid tumors	Phase I
ARCT-810	mRNA	Arcturus	Intravenously	Urea disorder	Phase I
CV7202	mRNA	CureVac	Intramuscular	Rabies	Phase I
mRNA-1893	mRNA	Moderna	Intramuscular	Zika	Phase I
CV9202	mRNA	CureVac	Intradermal	Non-small cell lung cancer	Phase I
mRNA-5671	mRNA	Moderna	Intravenous	Cancer	Phase I
BNT111	mRNA	BioNTech	Intravenous	Advanced Melanoma	Phase I

The table summarizes information available at ClinicalTrials.gov (<https://clinicaltrials.gov>) as of January 26, 2021. ASO, antisense oligonucleotide; siRNA, small interfering RNA; miRNA, microRNA; RNA, ribonucleic acid; mRNA, messenger RNA.



Gaviria, M., Kilic, B. A network analysis of COVID-19 mRNA vaccine patents. *Nat Biotechnol* **39**, 546–548 (2021). <https://doi.org/10.1038/s41587-021-00912-9>



01010...0010101010010  
001...01001010101011  
101...01001010101011  
010...01001010101010  
110...01001010101010  
10...01001010101011  
00...01001010101011  
01010101010101010  
11010101010101010010

COMPUTATIONAL  
AND STRUCTURAL  
BIOTECHNOLOGY  
JOURNAL

journal homepage: [www.elsevier.com/locate/csbj](http://www.elsevier.com/locate/csbj)



## TREAT: Therapeutic RNAs exploration inspired by artificial intelligence technology



Yufan Luo <sup>a,b,1</sup>, Liu Liu <sup>a,1</sup>, Zihao He <sup>a,1</sup>, Shanshan Zhang <sup>c</sup>, Peipei Huo <sup>c</sup>, Zhihao Wang <sup>c</sup>, Qin Jiaxin <sup>a</sup>, Lianhe Zhao <sup>a</sup>, Yang Wu <sup>a</sup>, Dongdong Zhang <sup>e</sup>, Dechao Bu <sup>a,d,\*</sup>, Runsheng Chen <sup>e,f,\*</sup>, Yi Zhao <sup>a,\*</sup>

<sup>a</sup> Research Center for Ubiquitous Computing Systems, Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China

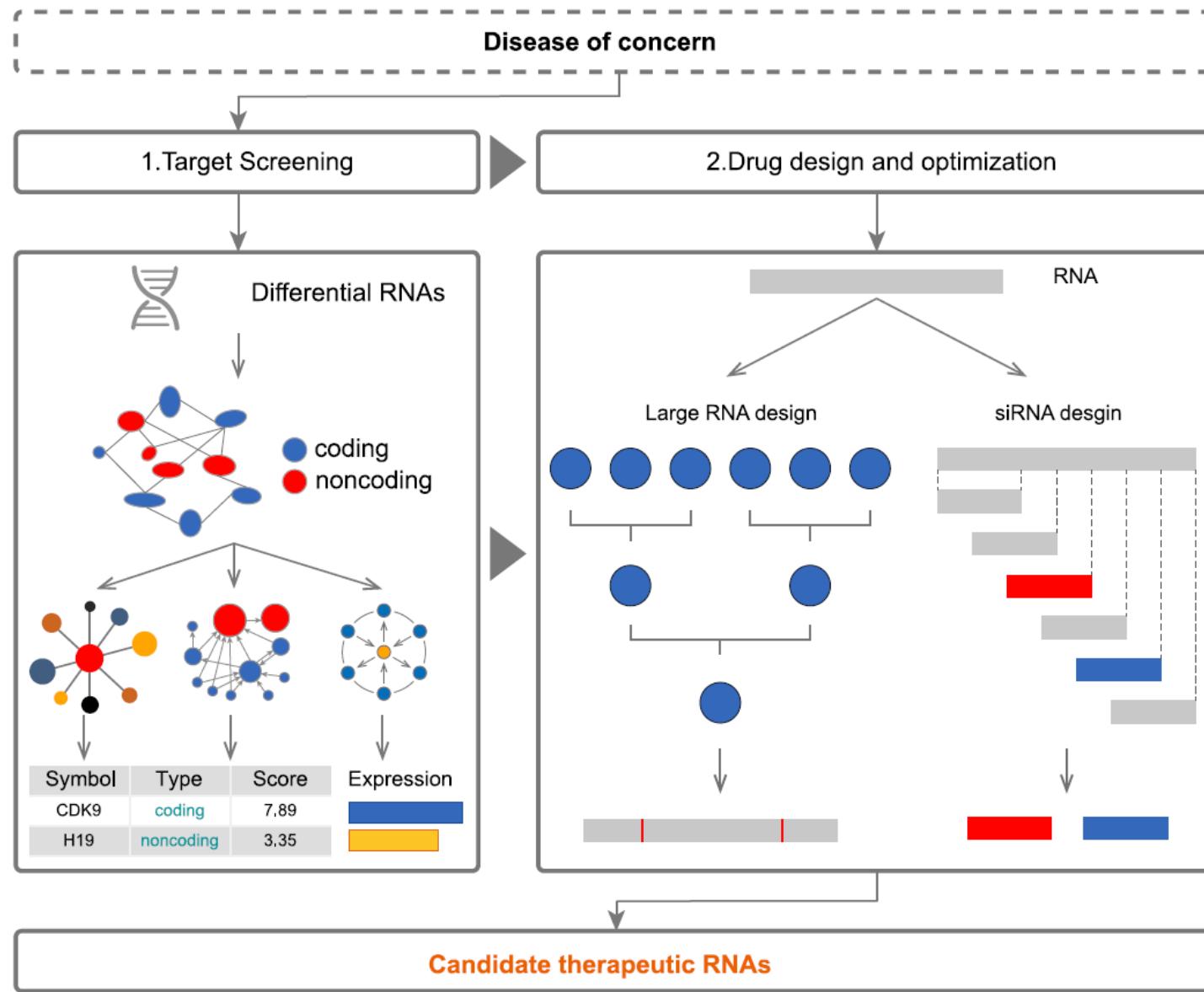
<sup>b</sup> University of Chinese Academy of Sciences, Beijing 100049, China

<sup>c</sup> Luoyang Zhongke Information Industry Research Institute, Luoyang, China

<sup>d</sup> Hwa Mei Hospital, University of Chinese Academy of Sciences, China

<sup>e</sup> Key Laboratory of RNA Biology, Center for Big Data Research in Health, Institute of Biophysics, Chinese Academy of Sciences, Beijing 100101, China

<sup>f</sup> Shenzhen Institute of Nucleic Acid Drug Research, Shenzhen Bay Laboratory Pingshan Translational Medicine Center, Shenzhen 510800, China



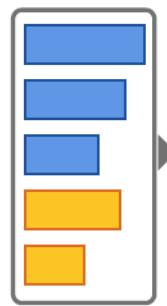
**A**

Input Seq

Bin 1                    Bin 2                    ...                    Bin n-1                    Bin n

ATGTTT	GTGTTT	...	CACTAC	ACATAA
--------	--------	-----	--------	--------

Features



Local  
Global

*Large RNA optimization*

Local optimization

Global optimization

Output Seq1

ATGTTT	GTGTC	...	CACTC	ACATAA
--------	-------	-----	-------	--------

Output Seq2

ATGTTT	GTGTTT	...	CACTCG	ACATAA
--------	--------	-----	--------	--------



Article

<https://doi.org/10.1038/s42256-022-00571-8>

# Deep learning models for predicting RNA degradation via dual crowdsourcing

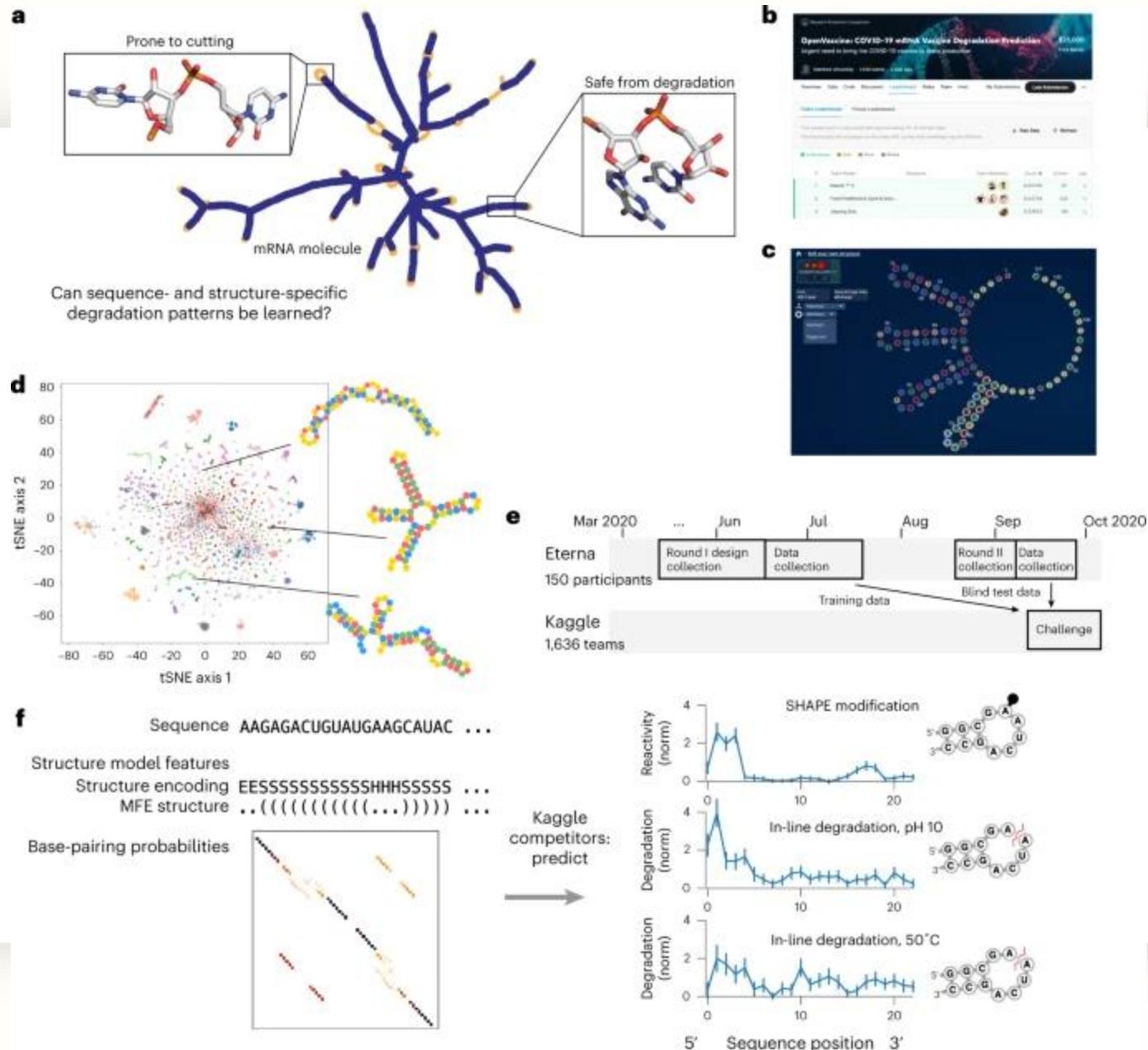
Received: 14 October 2021

Accepted: 21 October 2022

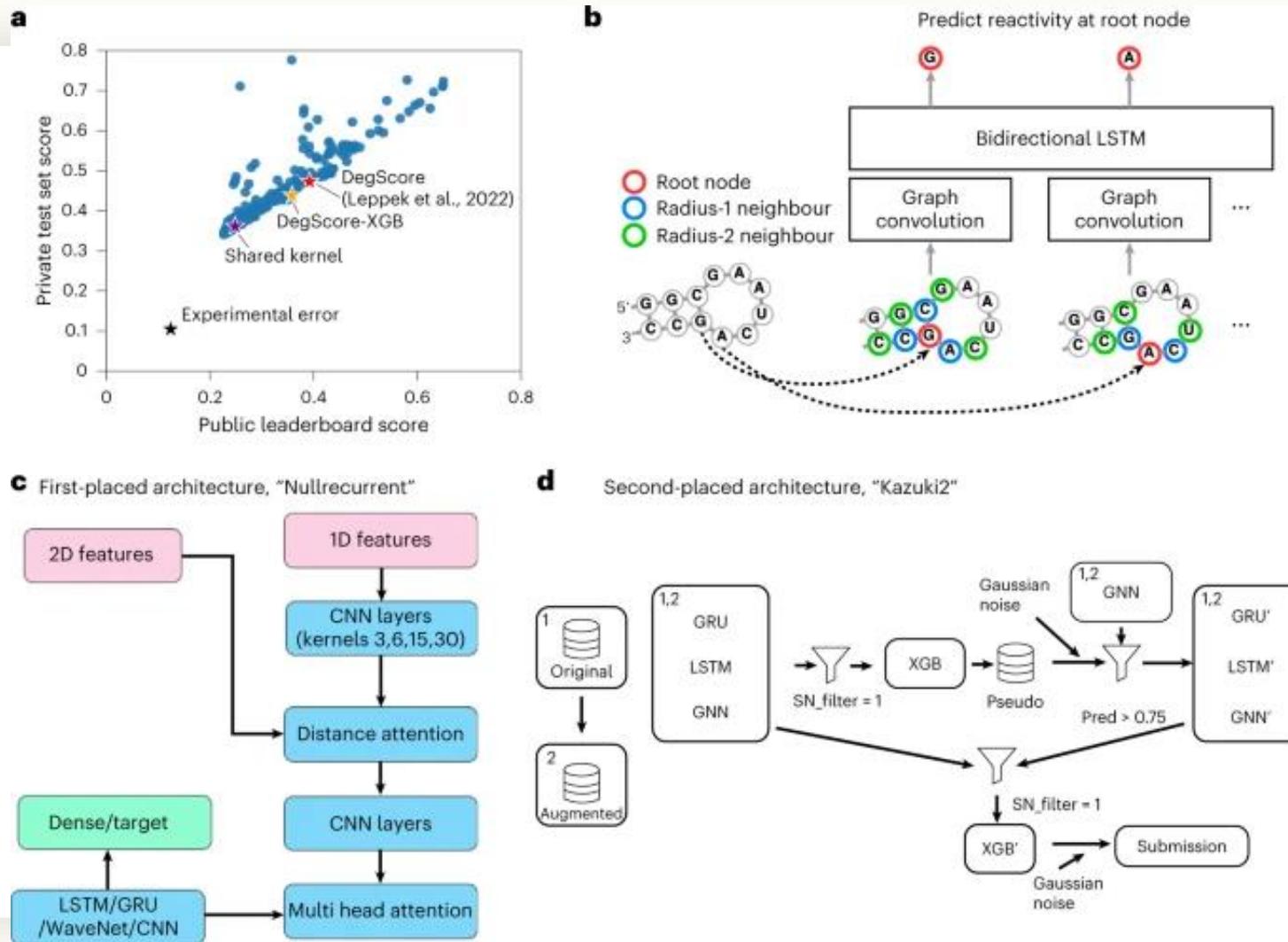
Published online: 14 December 2022

Check for updates

Hannah K. Wayment-Steele <sup>1,2,26</sup>, Wipapat Kladwang <sup>2,3,26</sup>,  
Andrew M. Watkins <sup>2,3,4,26</sup>, Do Soon Kim <sup>2,3,26</sup>, Bojan Tunguz <sup>3,5,26</sup>, Walter Reade <sup>6</sup>,  
Maggie Demkin <sup>6</sup>, Jonathan Romano <sup>2,3,7</sup>, Roger Wellington-Oguri <sup>2</sup>,  
John J. Nicol <sup>2</sup>, Jiayang Gao <sup>8</sup>, Kazuki Onodera <sup>9</sup>, Kazuki Fujikawa <sup>10</sup>,  
Hanfei Mao <sup>11</sup>, Gilles Vandewiele <sup>12</sup>, Michele Tinti <sup>13</sup>, Bram Steenwinckel <sup>12</sup>,  
Takuya Ito <sup>14</sup>, Taiga Noumi <sup>15</sup>, Shujun He <sup>16</sup>, Keiichiro Ishii <sup>17</sup>, Youhan Lee <sup>18,19</sup>,  
Fatih Öztürk <sup>20</sup>, King Yuen Chiu <sup>21</sup>, Emin Öztürk <sup>22</sup>, Karim Amer <sup>23</sup>,  
Mohamed Fares <sup>23,24</sup>, Eterna Participants\* & Rhiju Das <sup>2,3,25</sup>



\* Source: xxxx, 2022.0x.xx



\* Source: xxxx, 2022.0x.xx



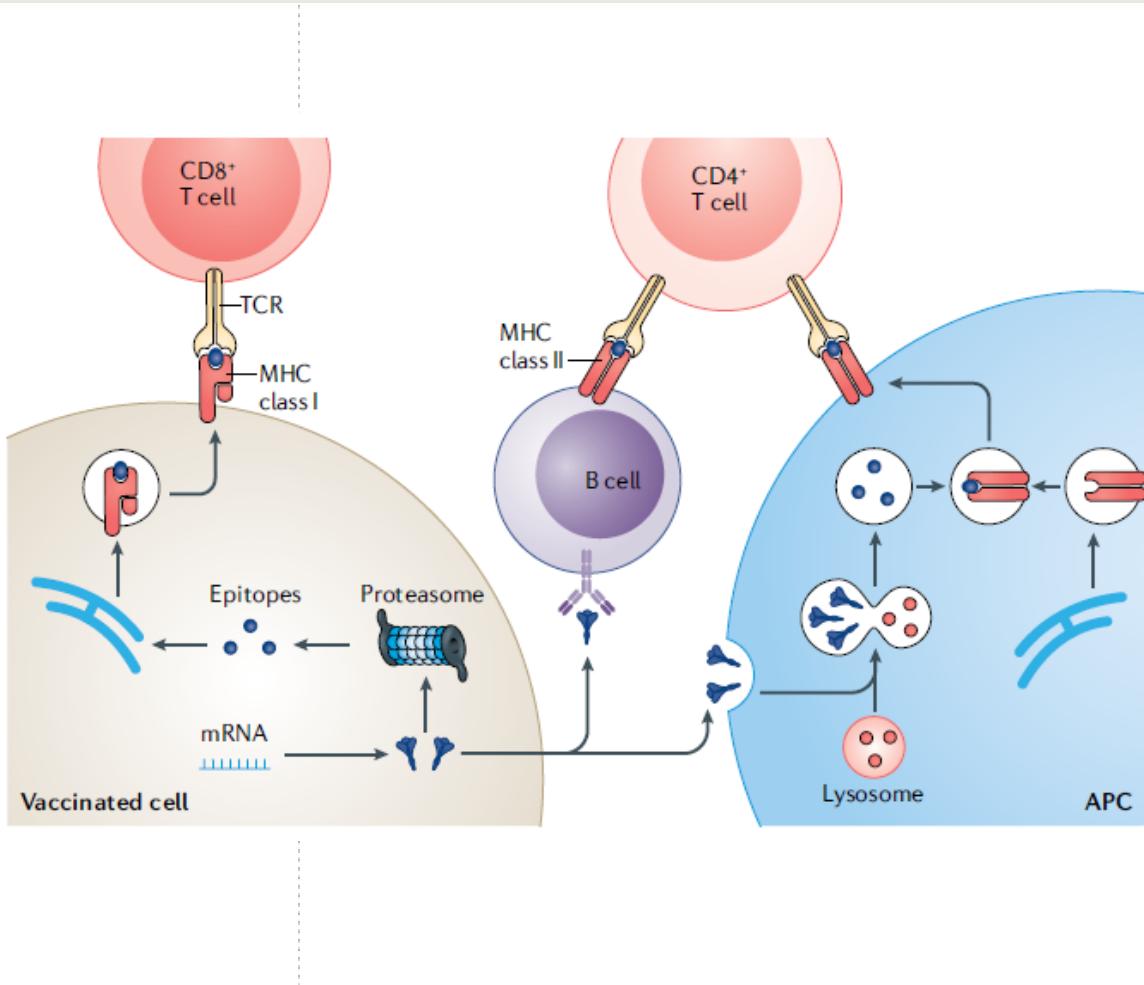
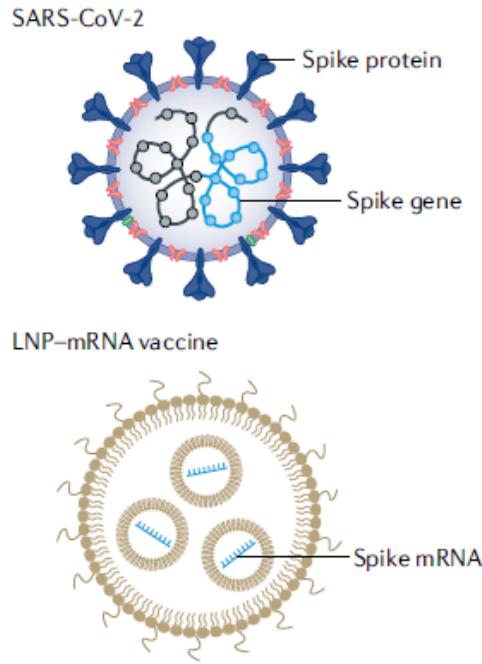
# REVIEWS

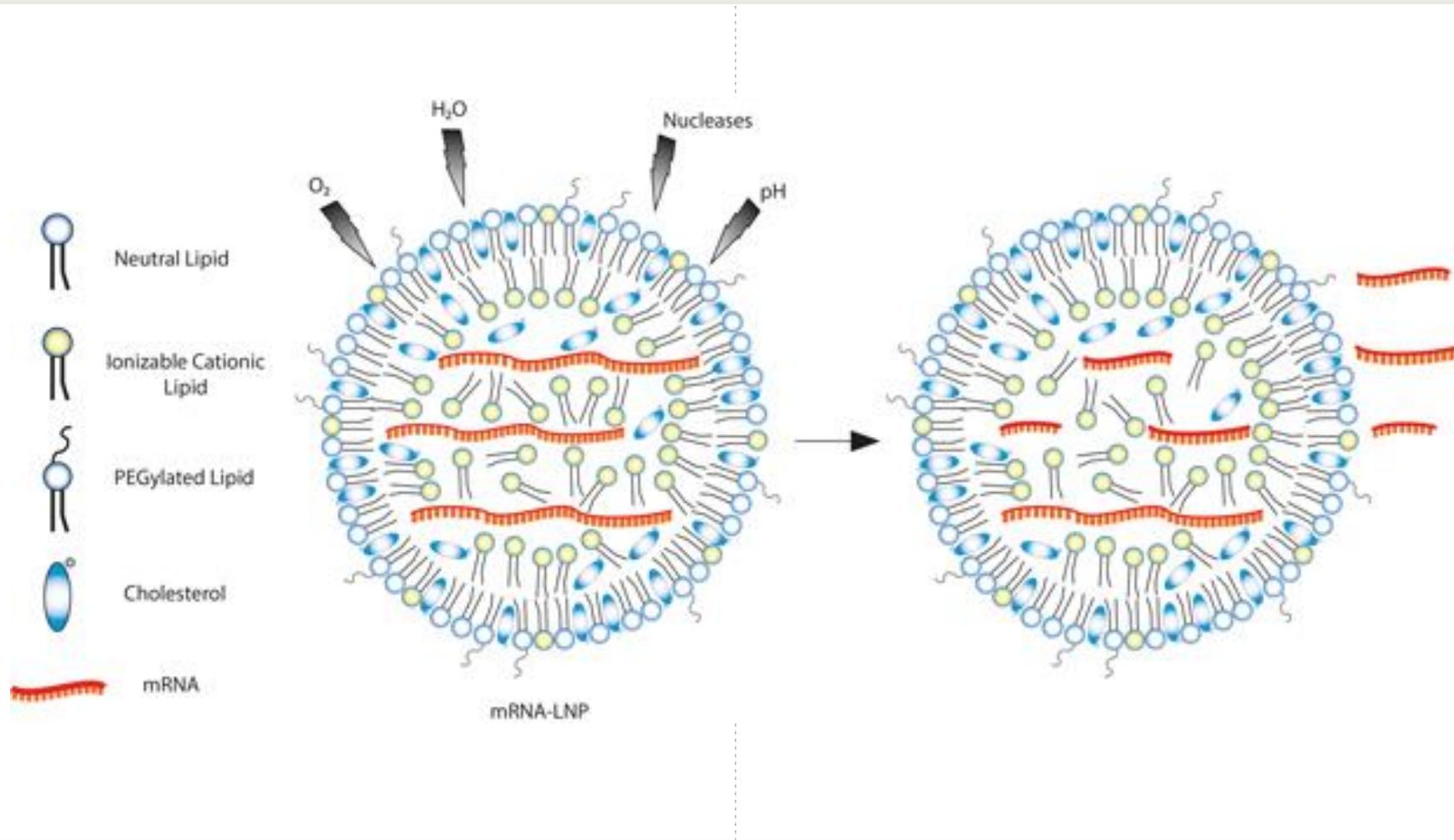
 Check for updates

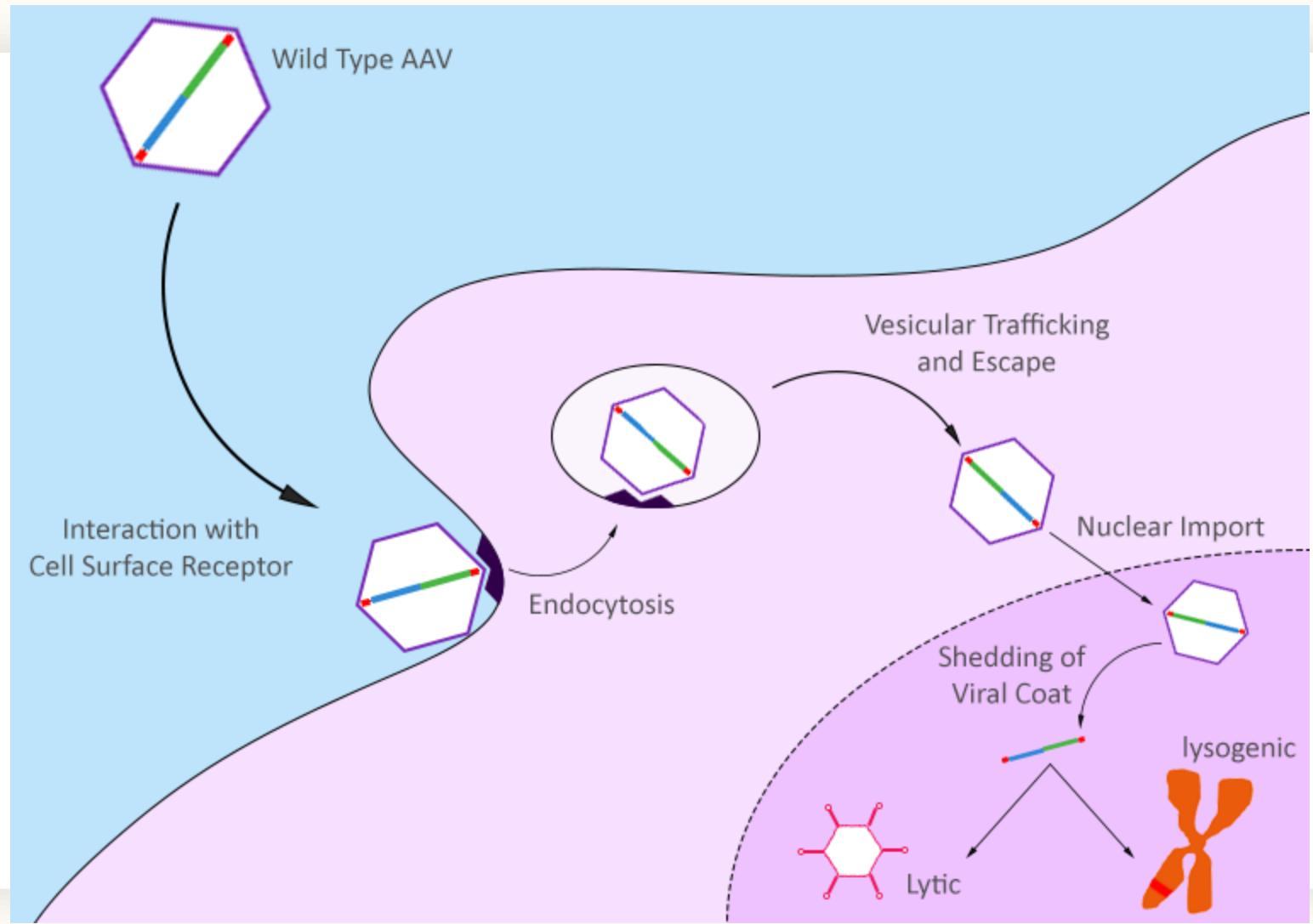
## Lipid nanoparticles for mRNA delivery

Xucheng Hou<sup>1</sup>, Tal Zaks<sup>2</sup> , Robert Langer<sup>1,3,4</sup>  and Yizhou Dong<sup>1</sup> 

**Abstract** | Messenger RNA (mRNA) has emerged as a new category of therapeutic agent to prevent and treat various diseases. To function *in vivo*, mRNA requires safe, effective and stable delivery systems that protect the nucleic acid from degradation and that allow cellular uptake and mRNA release. Lipid nanoparticles have successfully entered the clinic for the delivery of mRNA; in particular, lipid nanoparticle–mRNA vaccines are now in clinical use against coronavirus disease 2019 (COVID-19), which marks a milestone for mRNA therapeutics. In this Review, we discuss the design of lipid nanoparticles for mRNA delivery and examine physiological barriers and possible administration routes for lipid nanoparticle–mRNA systems. We then consider key points for the clinical translation of lipid nanoparticle–mRNA formulations, including good manufacturing practice, stability, storage and safety, and highlight preclinical and clinical studies of lipid nanoparticle–mRNA therapeutics for infectious diseases, cancer and genetic disorders. Finally, we give an outlook to future possibilities and remaining challenges for this promising technology.







# Deep-learning models for lipid nanoparticle-based drug delivery

Philip J Harrison<sup>\*,1</sup> , Håkan Wieslander<sup>2</sup> , Alan Sabirsh<sup>3</sup> , Johan Karlsson<sup>4</sup> , Victor Malmsjö<sup>1</sup> , Andreas Hellander<sup>2</sup> , Carolina Wählby<sup>2</sup>  & Ola Spjuth<sup>1</sup> 

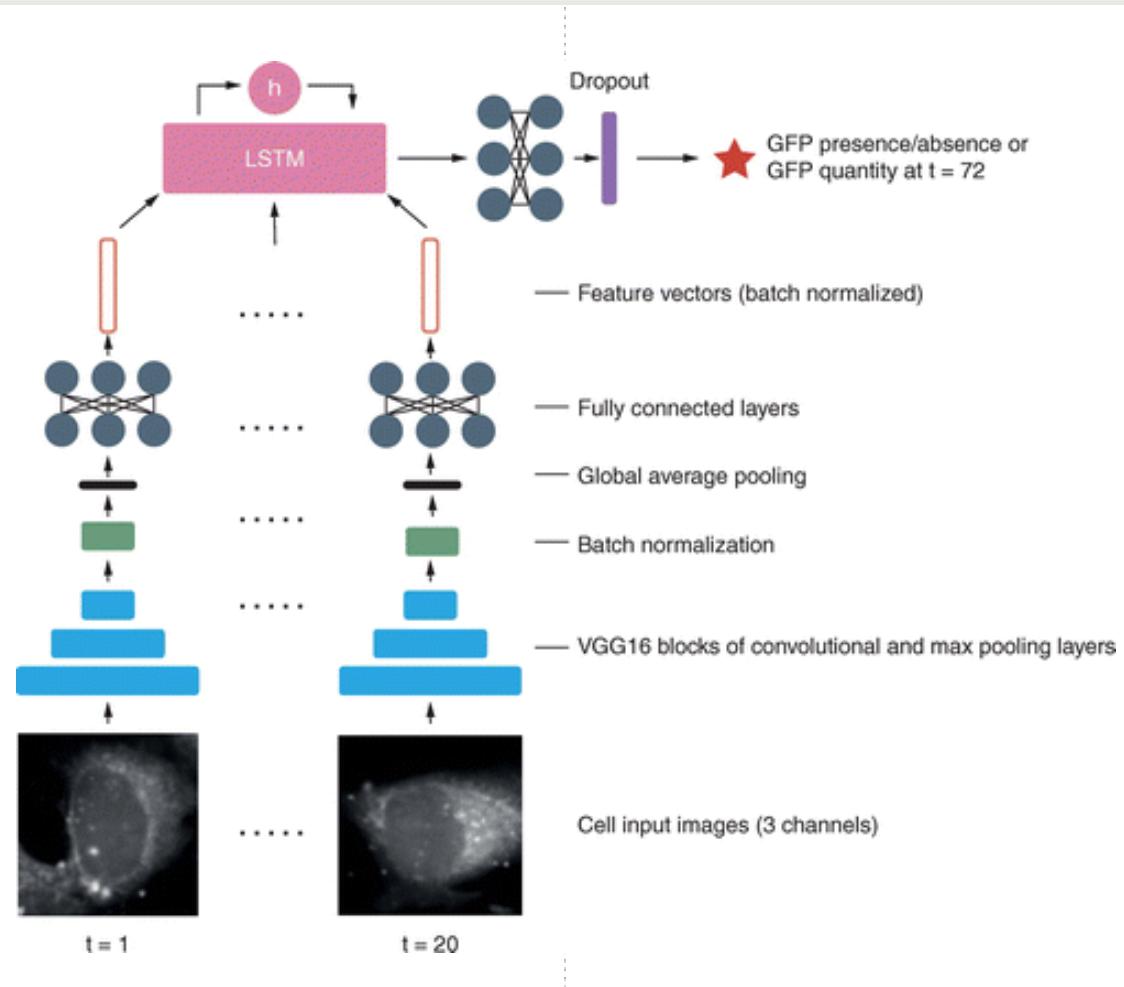
<sup>1</sup>Department of Pharmaceutical Biosciences, Uppsala University, Uppsala, Sweden

<sup>2</sup>Department of Information Technology, Uppsala University, Uppsala, Sweden

<sup>3</sup>Advanced Drug Delivery, Pharmaceutical Sciences, Research and Development, AstraZeneca, Gothenburg, Sweden

<sup>4</sup>Quantitative Biology, Discovery Sciences, Research and Development, AstraZeneca, Gothenburg, Sweden

\*Author for correspondence: philip.harrison@farmbio.uu.se

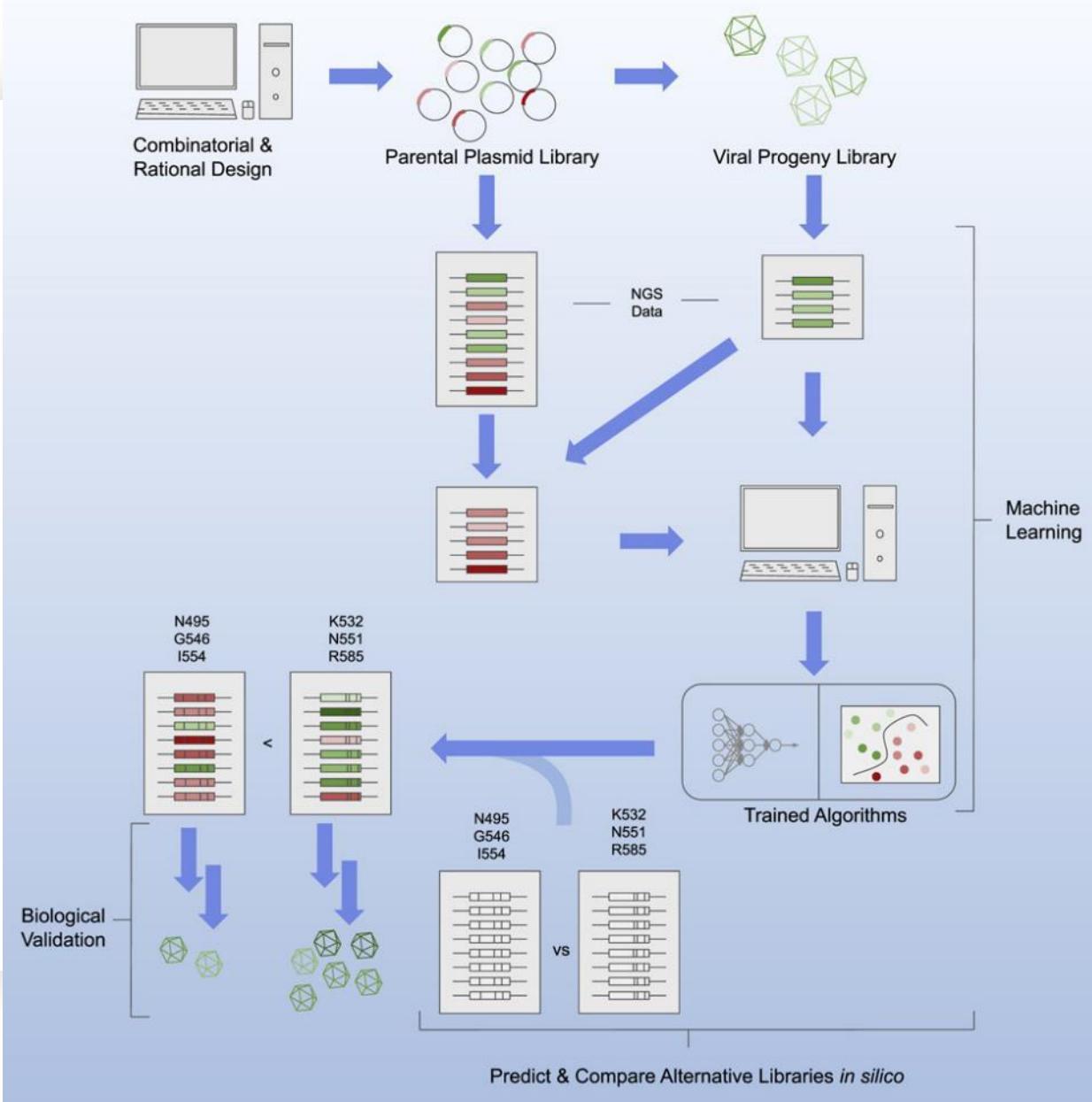




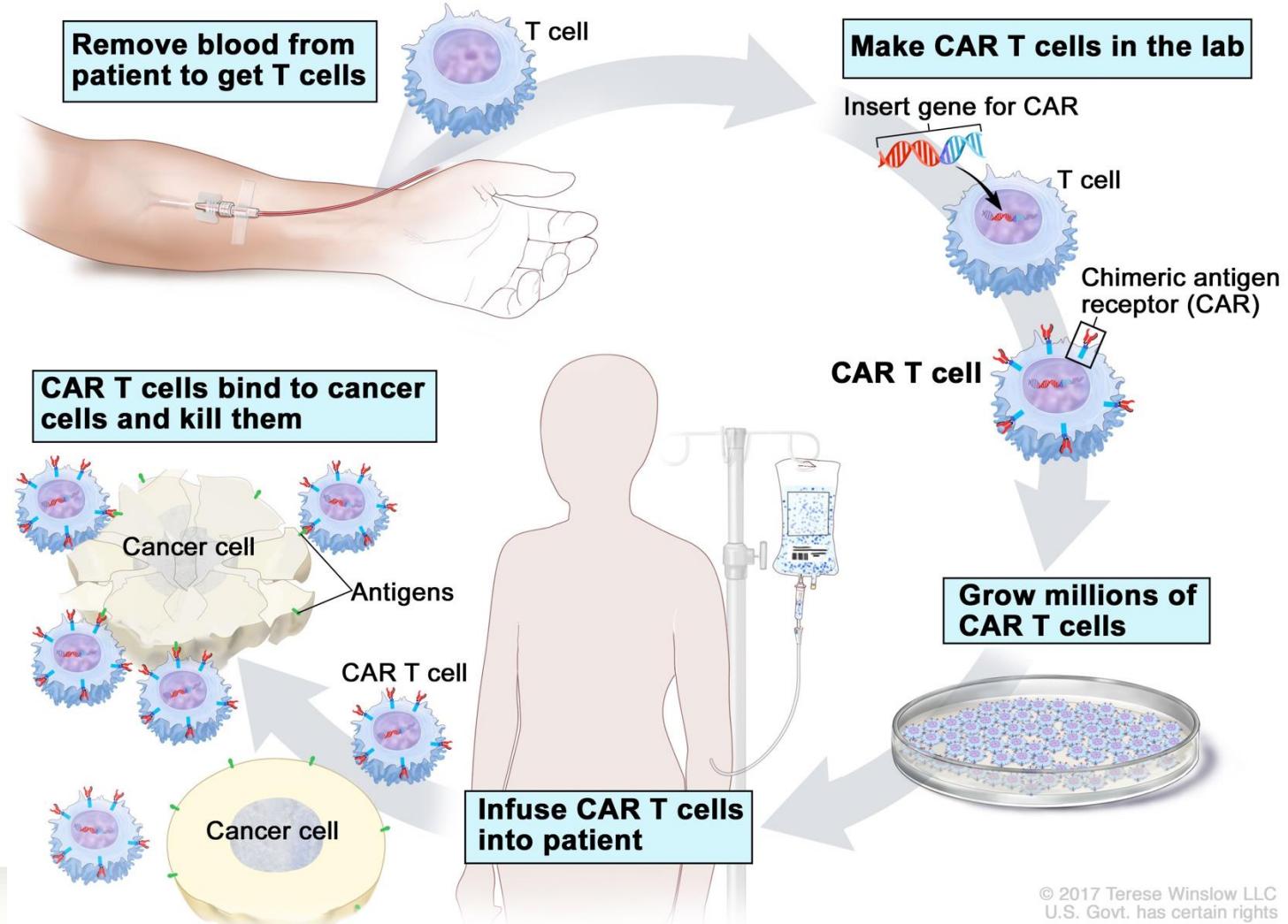
# Applying machine learning to predict viral assembly for adeno-associated virus capsid libraries

Andrew D. Marques,<sup>1</sup> Michael Kummer,<sup>2</sup> Oleksandr Kondratov,<sup>1</sup> Arunava Banerjee,<sup>2</sup> Oleksandr Moskalenko,<sup>3</sup>  
and Sergei Zolotukhin<sup>1</sup>

<sup>1</sup>Department of Pediatrics, Division of Cellular and Molecular Therapy, University of Florida, Gainesville, FL 32608, USA; <sup>2</sup>Department of Computer & Information Science & Engineering, University of Florida, Gainesville, FL 32603, USA; <sup>3</sup>University of Florida Research Computing, University of Florida, Gainesville, FL 32608, USA



## CAR T-cell Therapy

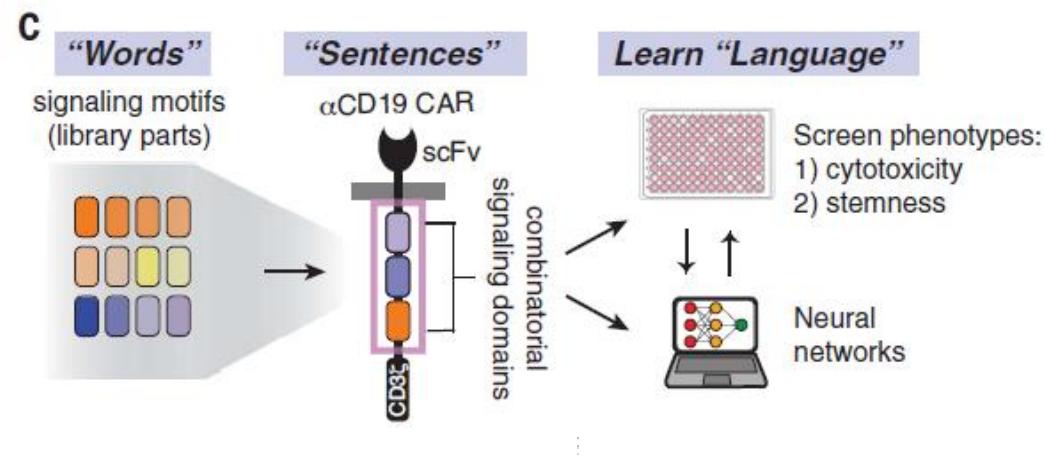
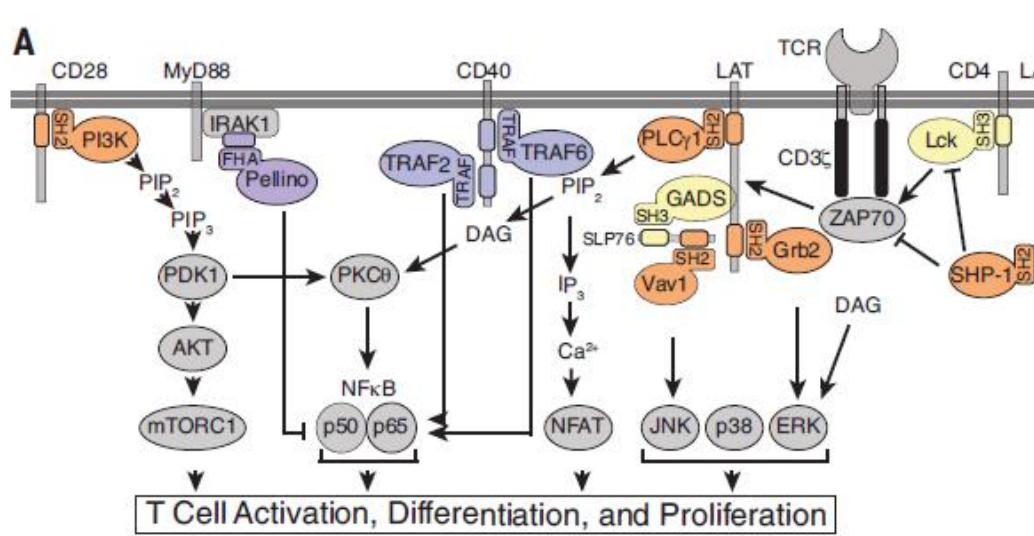


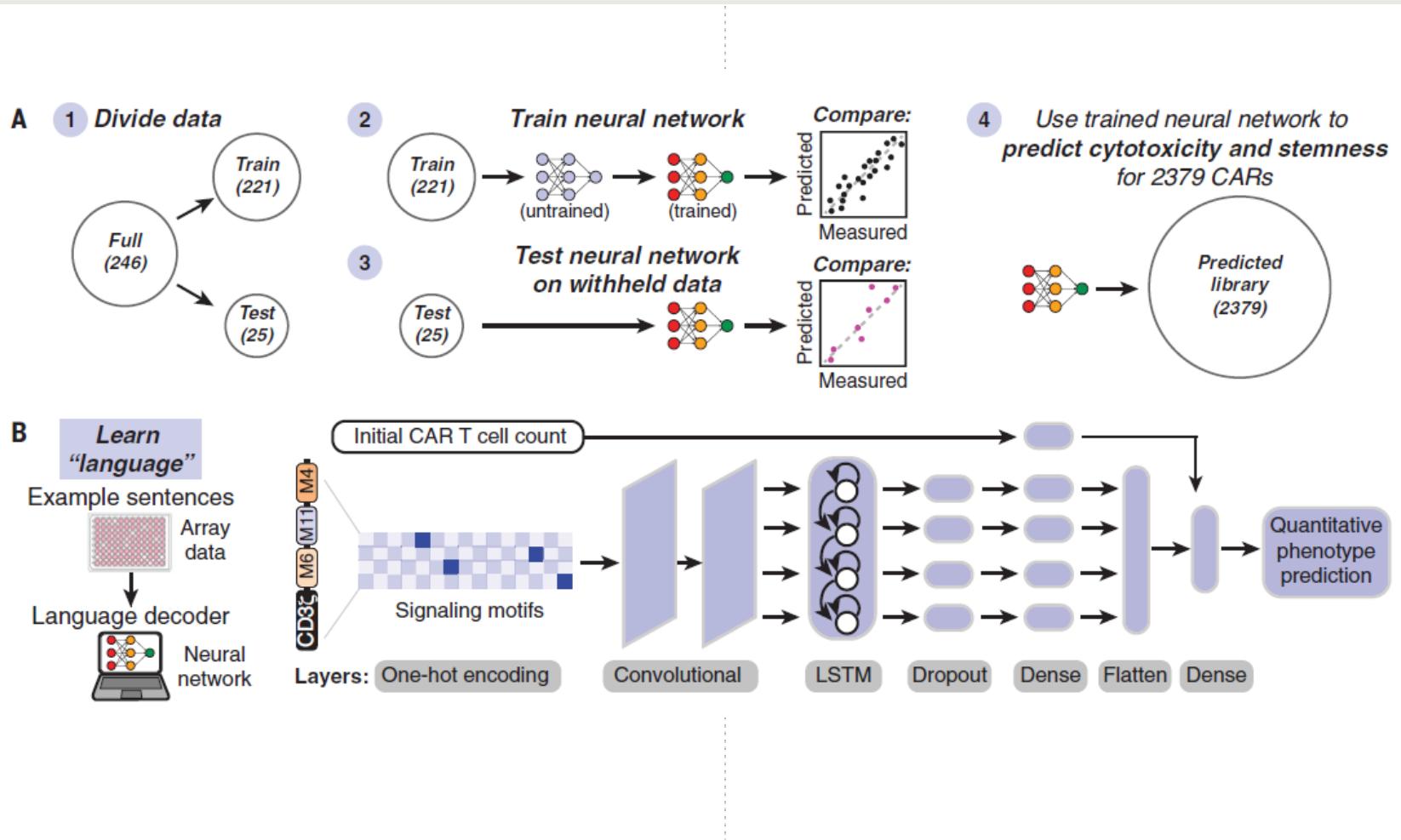
**SYNTHETIC BIOLOGY**

# Decoding CAR T cell phenotype using combinatorial signaling motif libraries and machine learning

Kyle G. Daniels<sup>1,2</sup>, Shangying Wang<sup>3,4</sup>, Milos S. Simic<sup>1,2</sup>, Hersh K. Bhargava<sup>1,2</sup>, Sara Capponi<sup>3,4</sup>, Yurie Tonai<sup>1,2</sup>, Wei Yu<sup>1,2</sup>, Simone Bianco<sup>3,4</sup>†\*, Wendell A. Lim<sup>1,2,4</sup>\*

Chimeric antigen receptor (CAR) costimulatory domains derived from native immune receptors steer the phenotypic output of therapeutic T cells. We constructed a library of CARs containing ~2300 synthetic costimulatory domains, built from combinations of 13 signaling motifs. These CARs promoted diverse human T cell fates, which were sensitive to motif combinations and configurations. Neural networks trained to decode the combinatorial grammar of CAR signaling motifs allowed extraction of key design rules. For example, non-native combinations of motifs that bind tumor necrosis factor receptor-associated factors (TRAFs) and phospholipase C gamma 1 (PLC $\gamma$ 1) enhanced cytotoxicity and stemness associated with effective tumor killing. Thus, libraries built from minimal building blocks of signaling, combined with machine learning, can efficiently guide engineering of receptors with desired phenotypes.



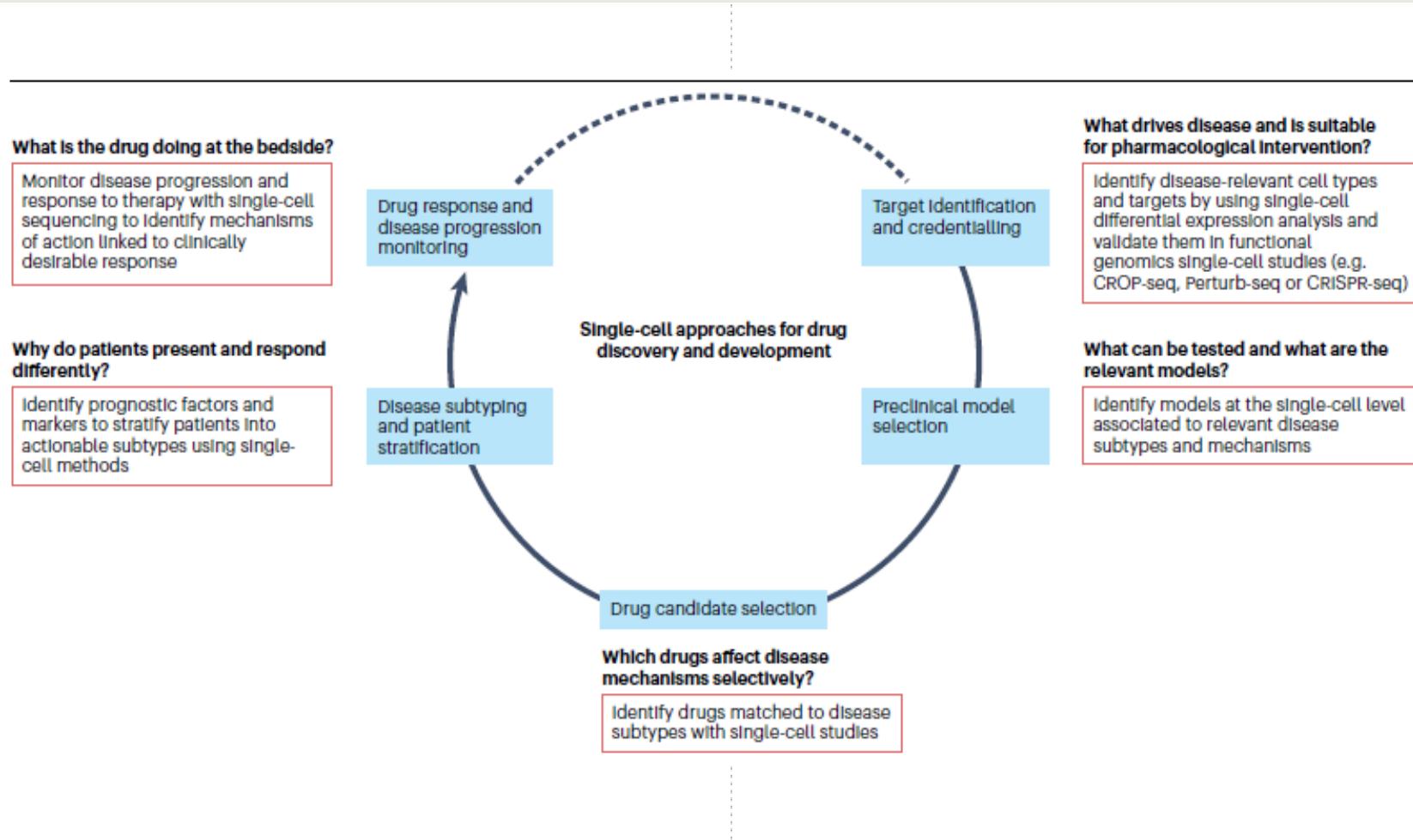


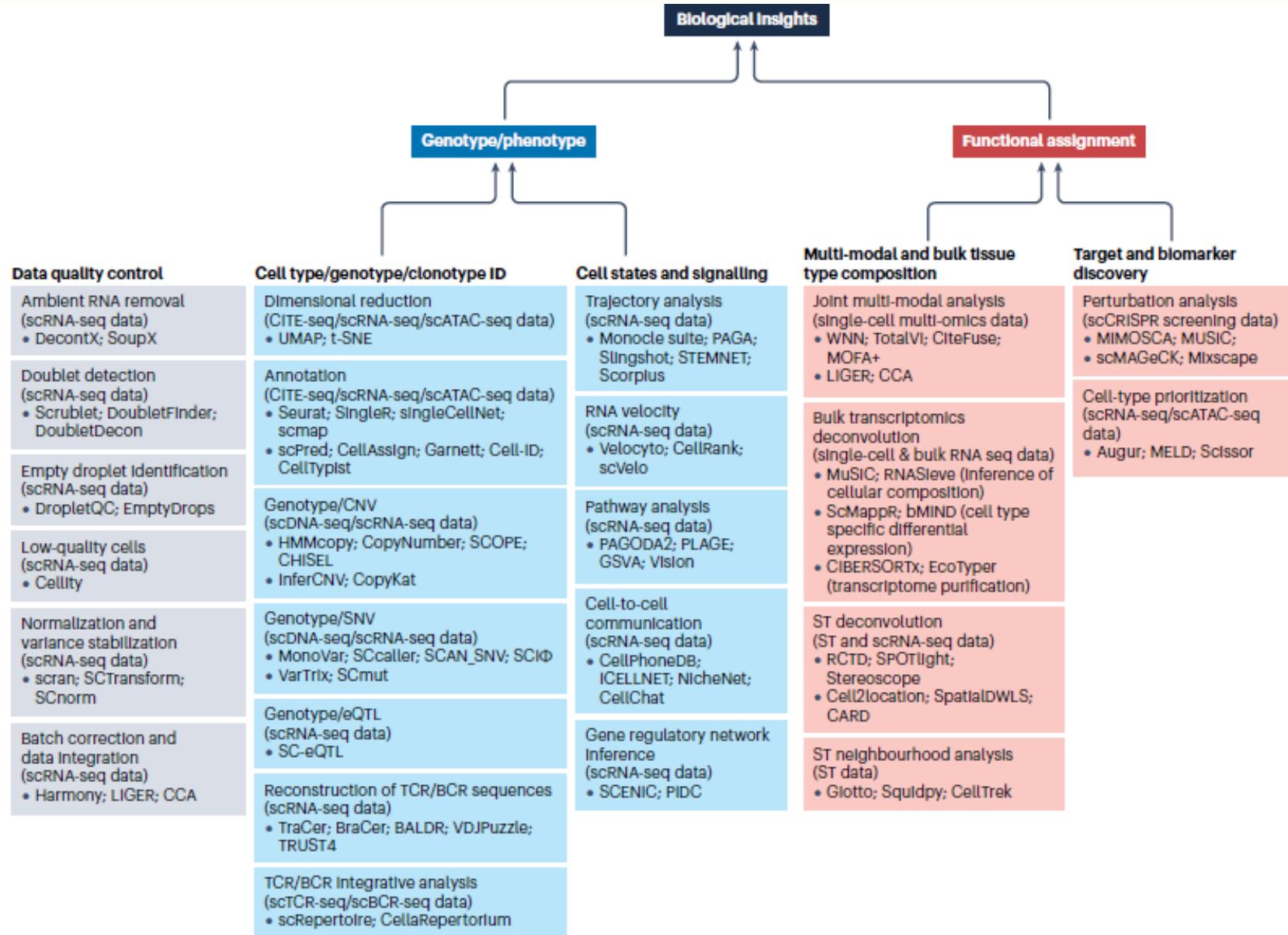
Review article

 Check for updates

# Applications of single-cell RNA sequencing in drug discovery and development

Bram Van de Sande<sup>1,15</sup>, Joon Sang Lee<sup>1,15</sup>, Euphemia Mutasa-Gottgens<sup>1,3,15</sup>✉, Bart Naughton<sup>1,4</sup>, Wendi Bacon<sup>1,3,5</sup>, Jonathan Manning<sup>3</sup>, Yong Wang<sup>1,6</sup>, Jack Pollard<sup>7</sup>, Melissa Mendez<sup>1,8</sup>, Jon Hill<sup>1,9</sup>, Namit Kumar<sup>1,10</sup>, Xiaohong Cao<sup>1,11</sup>, Xiao Chen<sup>12</sup>, Mugdha Khaladkar<sup>13</sup>, Ji Wen<sup>1,14</sup>, Andrew Leach<sup>1,3</sup> & Edgardo Ferran<sup>1,3</sup>





### a Standard high-throughput chemical screens

- 100k-1M compounds
- Typically one compound dose tested on one condition



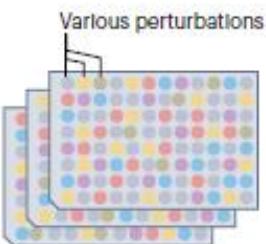
List of most active compounds

Dose-response and other studies

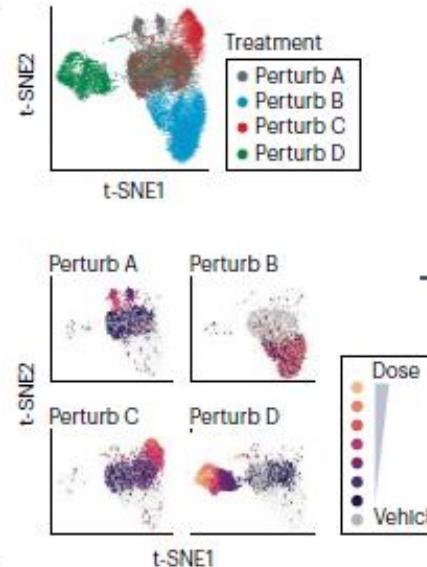
Hits suitable for drug discovery

### b Single-cell high-throughput chemical screens

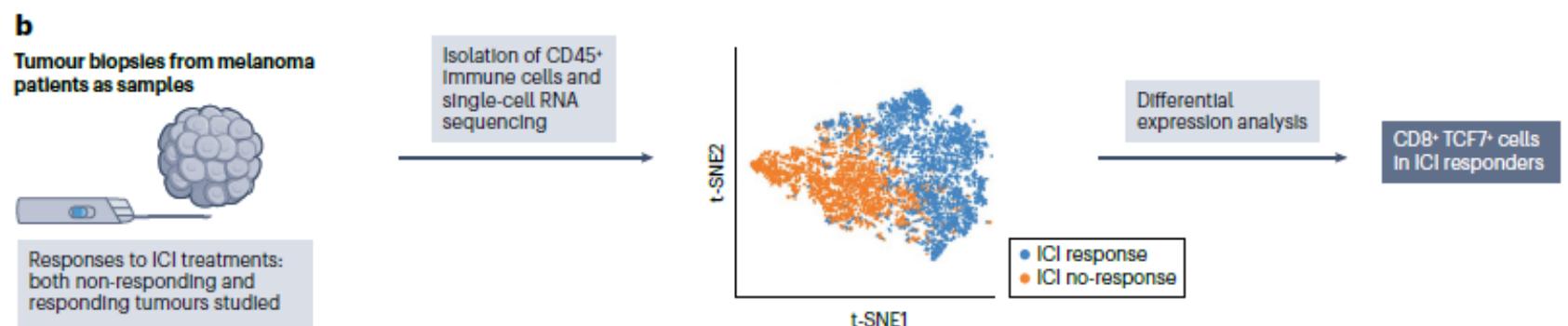
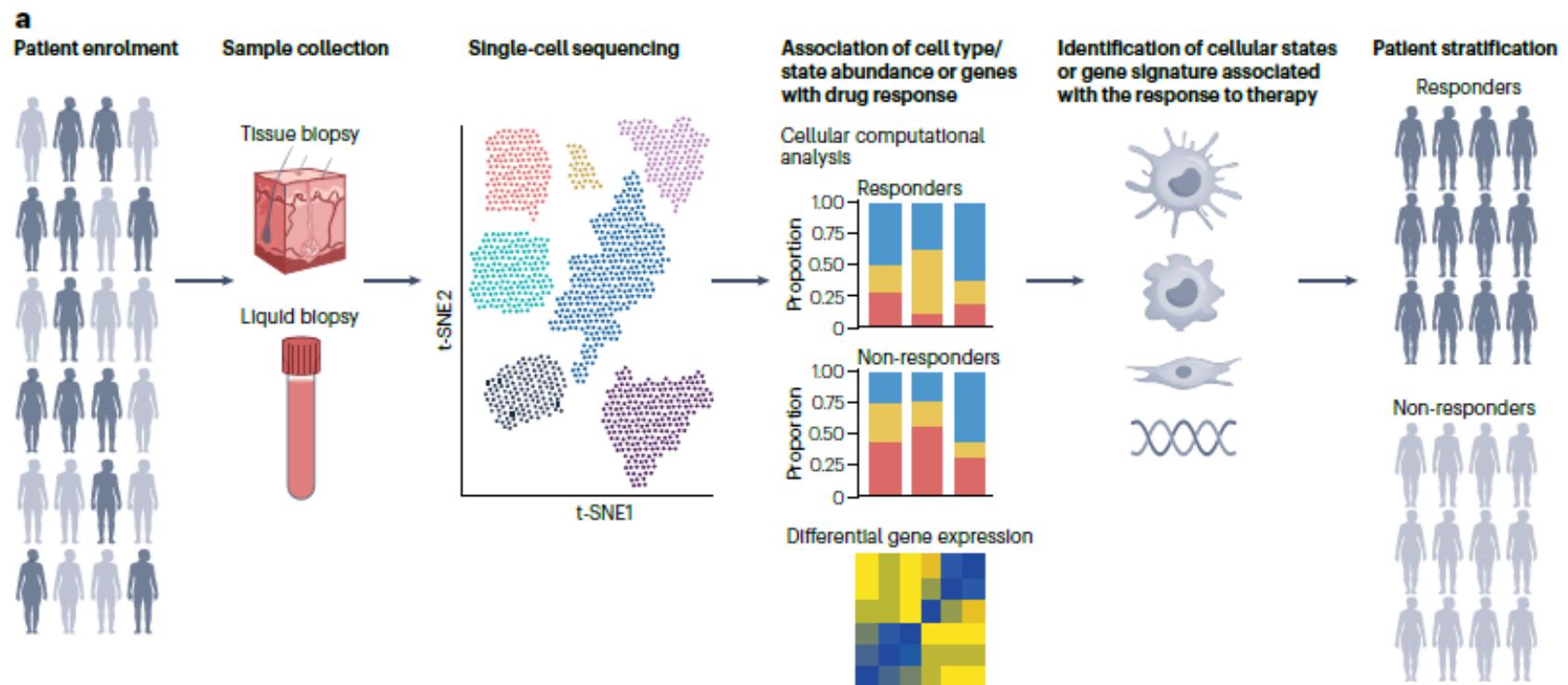
- 100–1,000 drugs
- Typically 5–10 drug doses tested on several cell lines (~100k-1M single cells in all)



Pooling scRNA-seq



Hints on drug's MoA





# HHS Public Access

Author manuscript

*Trends Pharmacol Sci.* Author manuscript; available in PMC 2021 December 01.

Published in final edited form as:

*Trends Pharmacol Sci.* 2020 December ; 41(12): 1050–1065. doi:10.1016/j.tips.2020.10.004.

## Single-cell Techniques and Deep Learning in Predicting Drug Response

Zhenyu Wu<sup>1,4</sup>, Patrick J. Lawrence<sup>1,4</sup>, Anjun Ma<sup>1</sup>, Jian Zhu<sup>2</sup>, Dong Xu<sup>3</sup>, Qin Ma<sup>1,\$</sup>

<sup>1</sup>Department of Biomedical Informatics, The Ohio State University, OH, 43235, USA

<sup>2</sup>Department of Pathology, The Ohio State University, OH, 43235, USA

<sup>3</sup>Department of Electrical Engineering and Computer Science, and Christopher S. Bond Life Sciences Center, University of Missouri, Columbia, MO 65211, USA

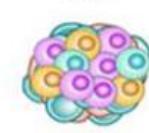
<sup>4</sup>These authors contributed equally to the paper as first authors

### Tumor profiling techniques

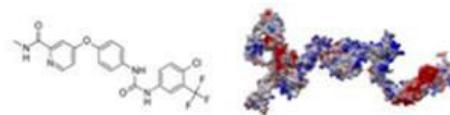
Single-cell



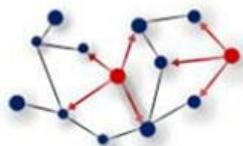
Bulk



### Drug chemical and structural data

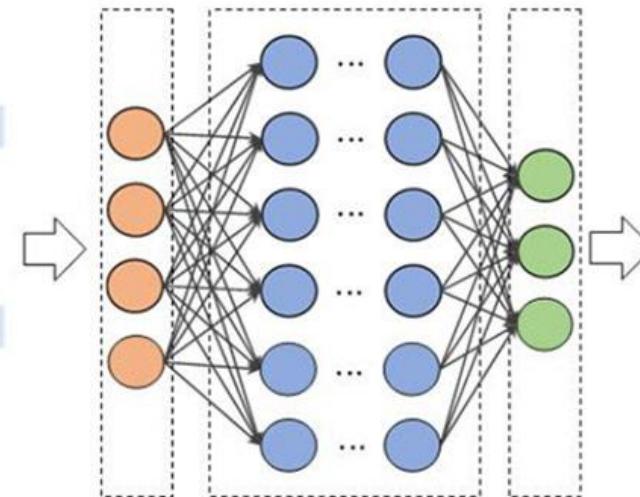


### Drug-Target (Ligand) Information



### Deep learning architecture

Input layer      Hidden layers      Output layer



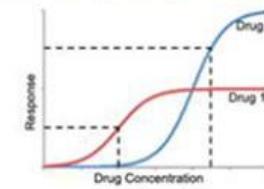
### Drug-target interactions

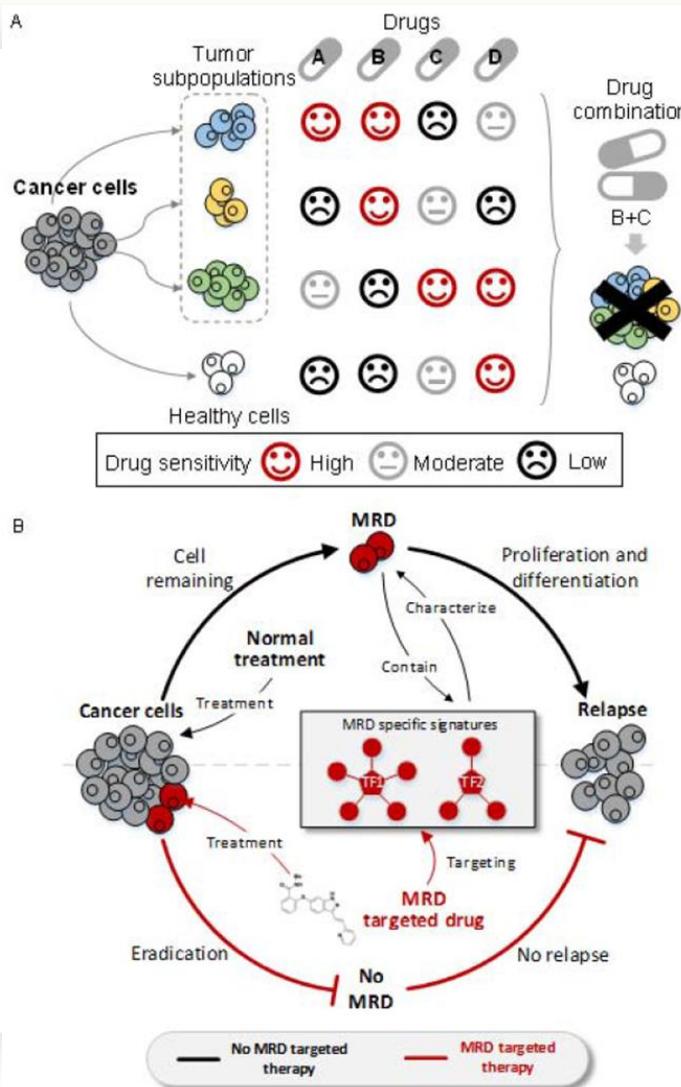


### Targeted drugs



### Drug response







# Deep transfer learning of cancer drug responses by integrating bulk and single-cell RNA-seq data

---

Received: 6 August 2021

Junyi Chen<sup>1,6</sup>, Xiaoying Wang<sup>2,6</sup>, Anjun Ma<sup>3,5</sup>✉, Qi-En Wang<sup>4</sup>, Bingqiang Liu<sup>2</sup>, Lang Li<sup>1</sup>, Dong Xu<sup>5</sup> & Qin Ma<sup>1,3</sup>✉

---

Accepted: 19 October 2022

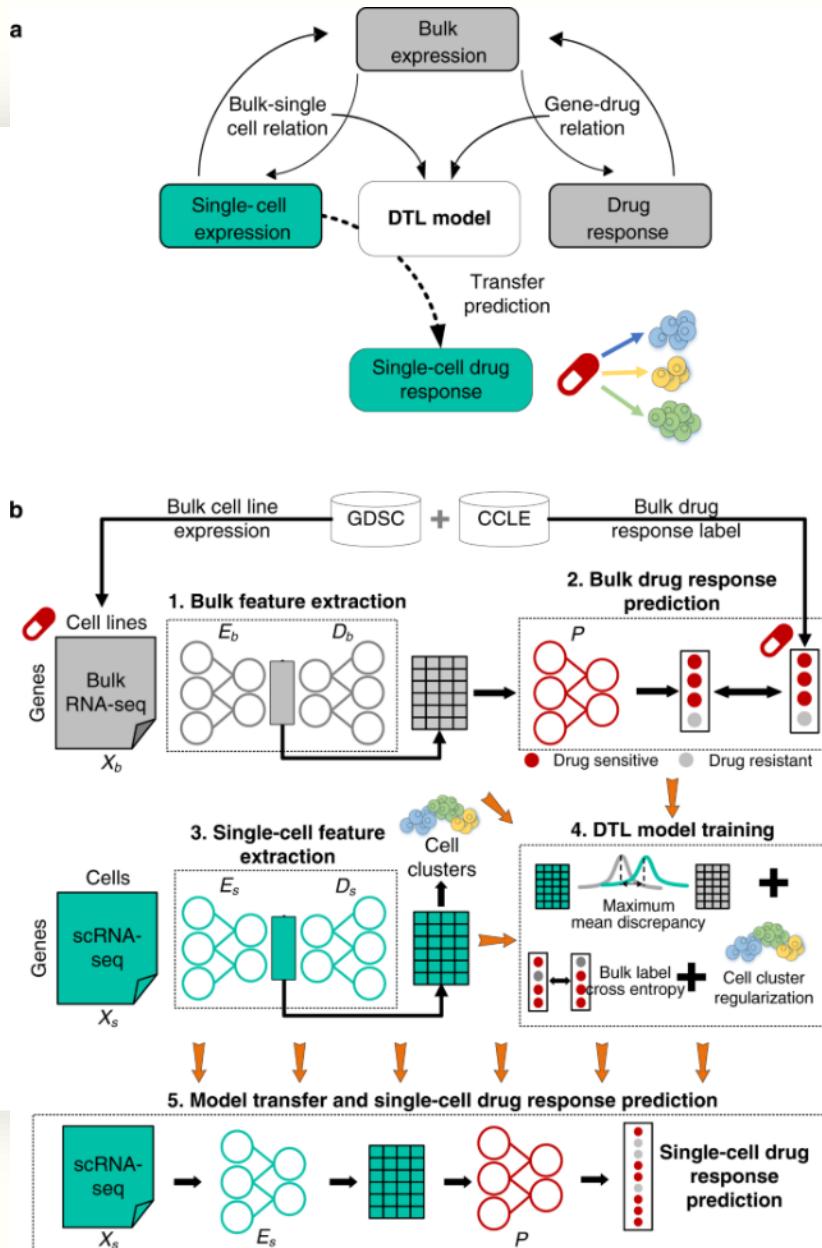
---

Published online: 30 October 2022

---

Check for updates

Drug screening data from massive bulk gene expression databases can be analyzed to determine the optimal clinical application of cancer drugs. The growing amount of single-cell RNA sequencing (scRNA-seq) data also provides insights into improving therapeutic effectiveness by helping to study the heterogeneity of drug responses for cancer cell subpopulations. Developing computational approaches to predict and interpret cancer drug response in single-cell data collected from clinical samples can be very useful. We propose scDEAL, a deep transfer learning framework for cancer drug response prediction at the single-cell level by integrating large-scale bulk cell-line data. The highlight in scDEAL involves harmonizing drug-related bulk RNA-seq data with scRNA-seq data and transferring the model trained on bulk RNA-seq data to predict drug responses in scRNA-seq. Another feature of scDEAL is the integrated gradient feature interpretation to infer the signature genes of drug resistance mechanisms. We benchmark scDEAL on six scRNA-seq datasets and demonstrate its model interpretability via three case studies focusing on drug response label prediction, gene signature identification, and pseudotime analysis. We believe that scDEAL could help study cell reprogramming, drug selection, and repurposing for improving therapeutic efficacy.



## REVIEW

# Trends and Potential of Machine Learning and Deep Learning in Drug Study at Single-Cell Level

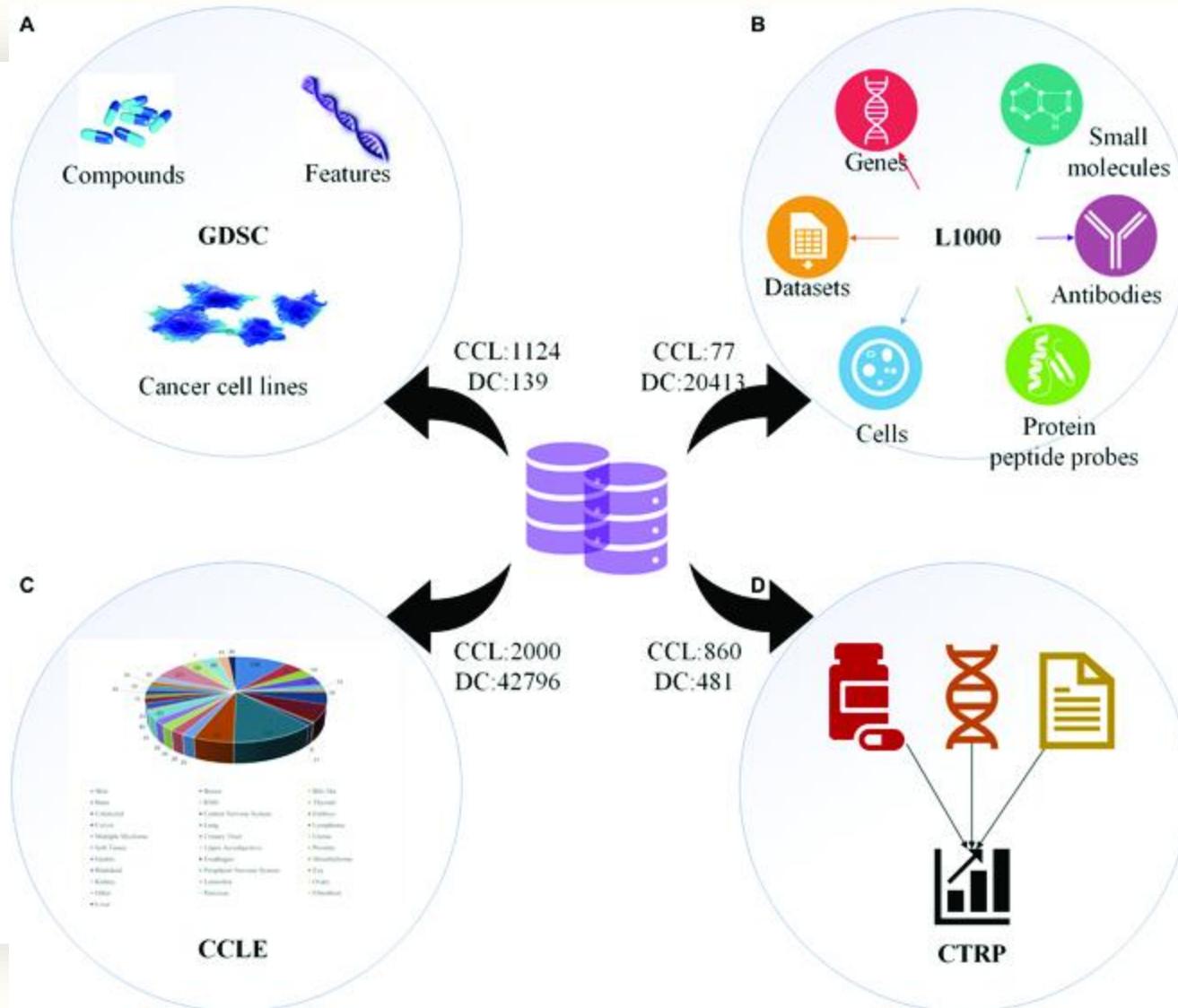
Ren Qi<sup>1,2</sup> and Quan Zou<sup>1,3\*</sup>

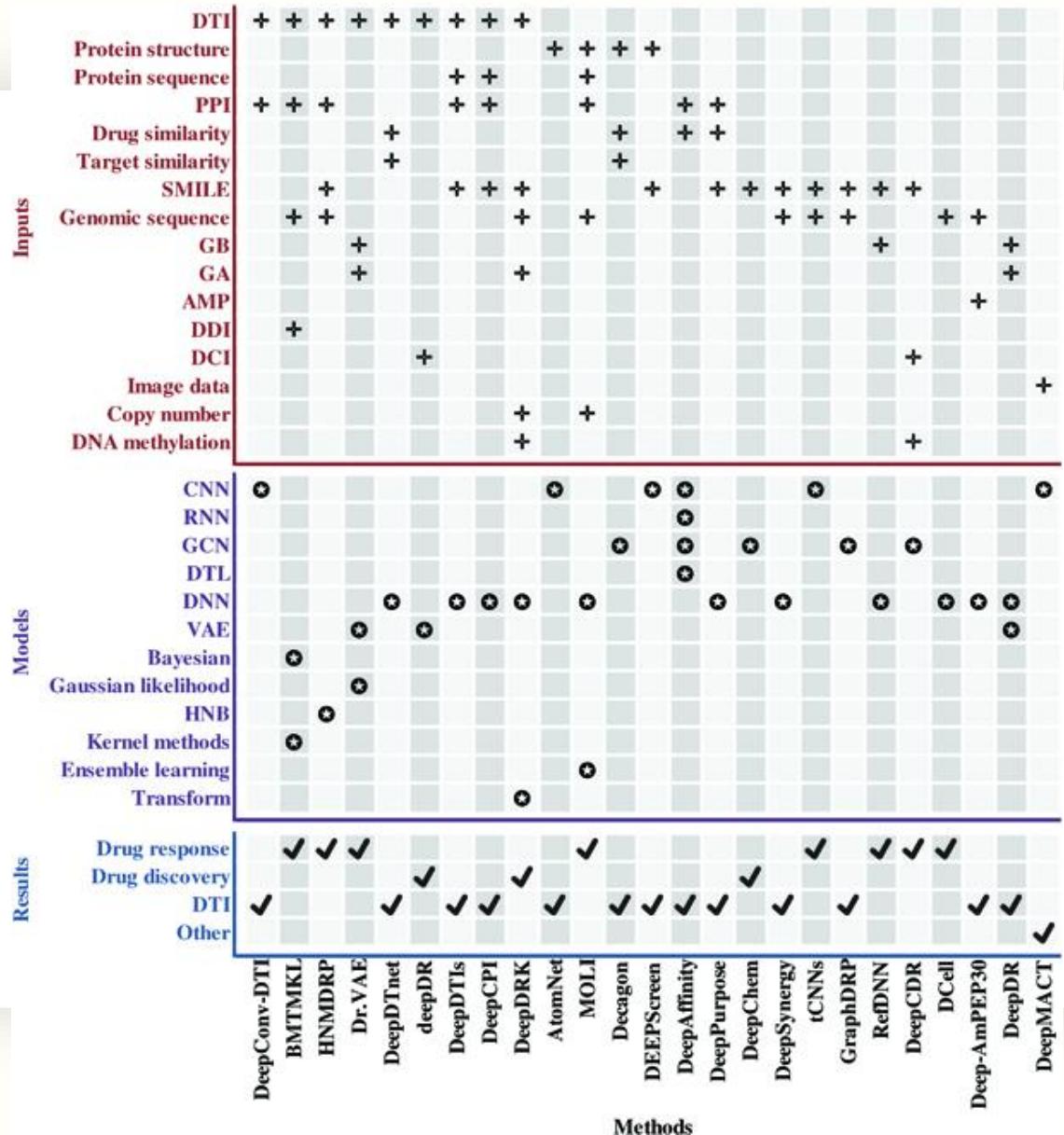
<sup>1</sup>Yangtze Delta Region Institute (Quzhou), University of Electronic Science and Technology of China, Quzhou 324000, China. <sup>2</sup>School of Life Science and Technology, University of Electronic Science and Technology of China, Chengdu, China. <sup>3</sup>Institute of Fundamental and Frontier Sciences, University of Electronic Science and Technology of China, Chengdu, China.

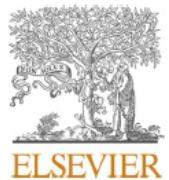
\*Address correspondence to: [zouquan@nclab.net](mailto:zouquan@nclab.net)

**Citation:** Qi R, Zou Q. Trends And Potential Of Machine Learning And Deep Learning In Drug Study At Single-Cell Level. *Research* 2023;6:Article 0050. <https://doi.org/10.34133/research.0050>

Submitted 22 September 2022  
Accepted 27 December 2022  
Published 9 March 2023







## Artificial intelligence in cancer immunotherapy: Applications in neoantigen recognition, antibody design and immunotherapy response prediction

Tong Li <sup>a,1</sup>, Yupeng Li <sup>a,1</sup>, Xiaoyi Zhu <sup>b,d,1</sup>, Yao He <sup>a</sup>, Yanling Wu <sup>b,d</sup>, Tianlei Ying <sup>b,d,\*†,1</sup>,  
Zhi Xie <sup>a,c,\*\*</sup>

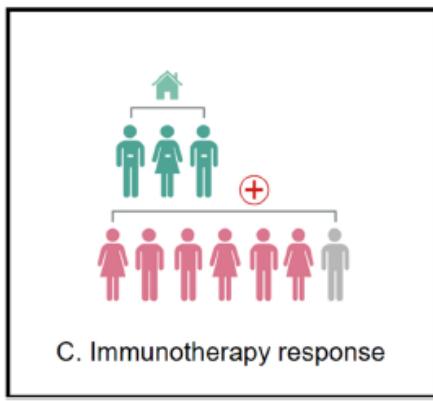
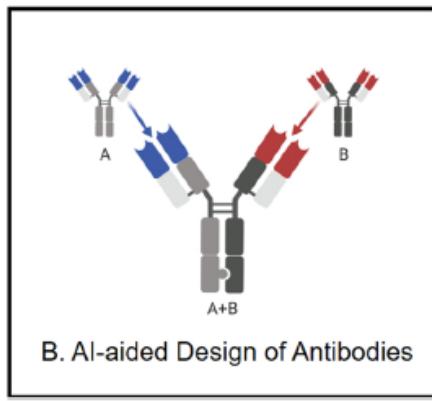
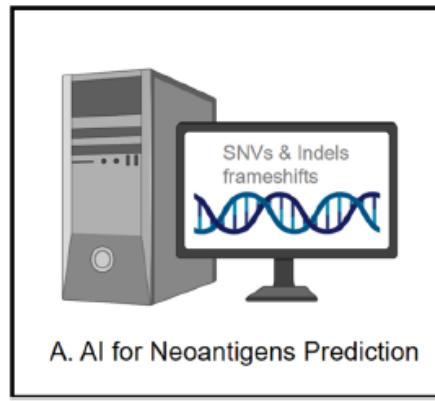
<sup>a</sup> State Key Laboratory of Ophthalmology, Zhongshan Ophthalmic Center, Sun Yat-sen University, Guangzhou, China

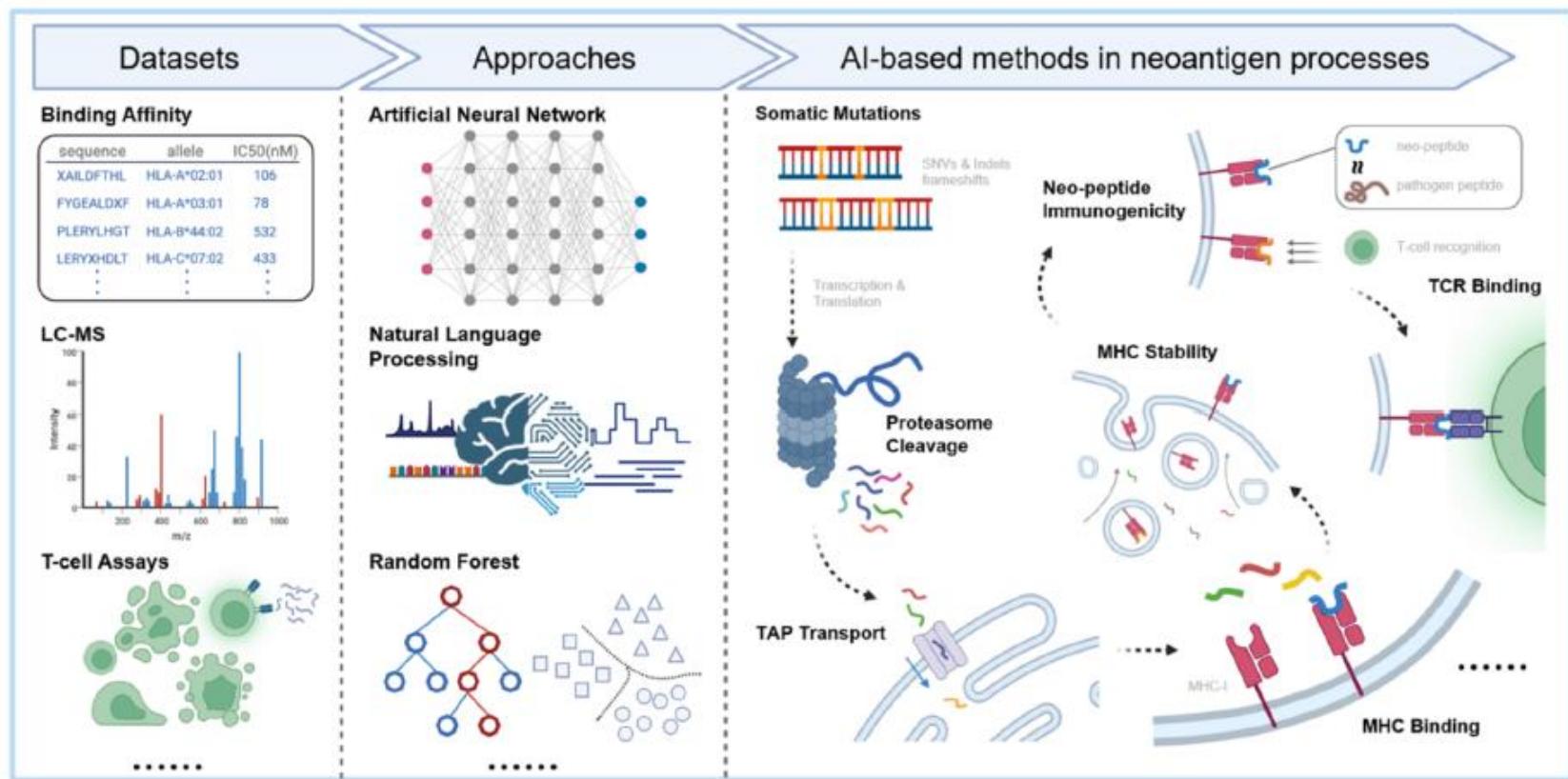
<sup>b</sup> MOE/NHC Key Laboratory of Medical Molecular Virology, Shanghai Institute of Infectious Disease and Biosecurity, School of Basic Medical Sciences, Shanghai Medical College, Fudan University, Shanghai, China

<sup>c</sup> Center for Precision Medicine, Sun Yat-sen University, Guangzhou, China

<sup>d</sup> Shanghai Engineering Research Center for Synthetic Immunology, Shanghai, China



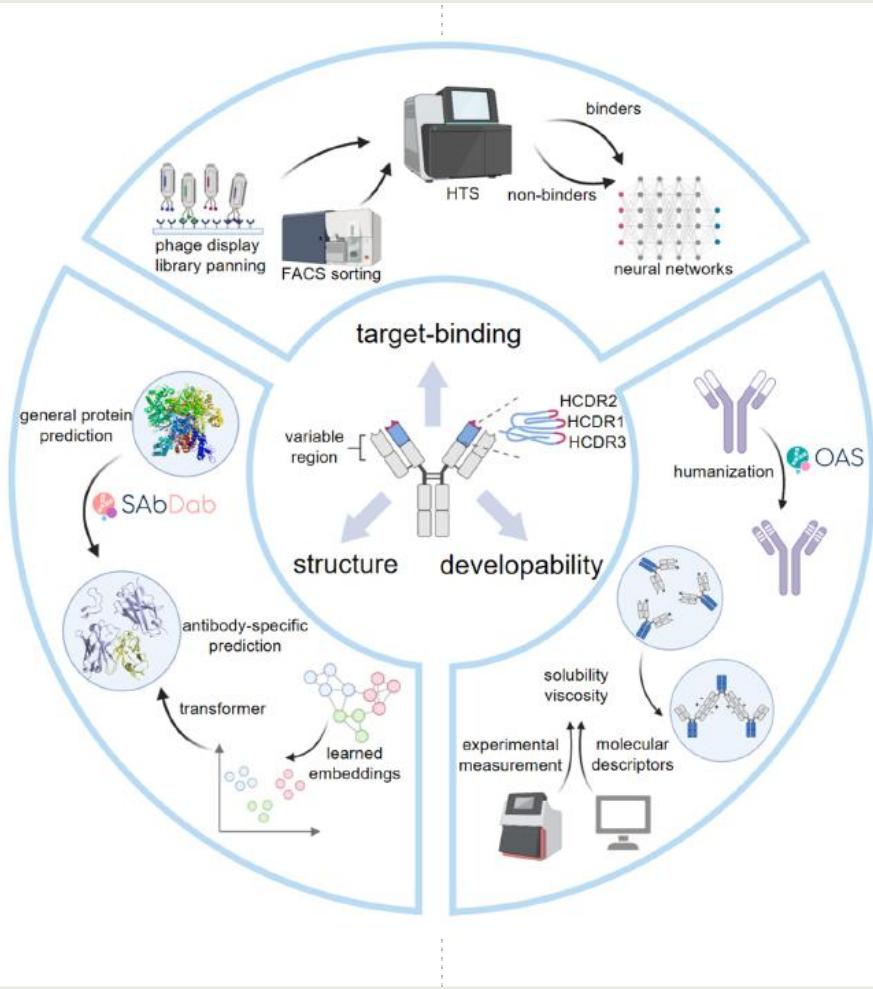




**Table 1**

Application of different AI algorithms in immunotherapy.

1. AI-based methods for neoantigen prediction								
Name	Year	MHC	Approach	Neoantigen process	Type of learning data	Training sets	Independent benchmark	
<i>EDGE</i> [38]	2019	I/II	DL	MHC binding	MHC ligands	142,844 peptides from 101 samples		
<i>DeepHLApan</i> [225]	2019	I	RNN	MHC binding & TCR binding	MHC ligands & T-cell recognized neoepitopes	437,077 peptides for binding model & 32,785 peptides for immunogenic		
<i>NMMP</i> [63]	2021	I	RF	TCR binding	T-cell recognized neoepitopes	185 neoepitopes identified by TILs from 96 samples		
<i>NetMHC-4.0</i> [47]	2016	I	NN	MHC binding	Affinity measurements	Affinity measurements from IEDB	[44,46]	
<i>NetMHCpan-4.0</i> [36]	2017	I	NN	MHC binding	Affinity measurements & MHC ligands	Affinity measurements from IEDB & 85,217 peptides	[44,46]	
<i>NetMHCstabpan</i> [67]	2016	I	NN	pMHC stability	Half-life of pMHC in vitro	28,166 binding measurements	[44]	
<i>MHCflurry</i> [50]	2018	I	NN	MHC binding	Affinity measurements & MHC ligands	230,735 affinity measurements & 226,684 peptides	[44,46]	
<i>MHCflurry-2.0</i> [51]	2020	I	NN	binding affinity & MHC ligands	Affinity measurements & MHC ligands	219,596 affinity measurements & 493,473 peptides for binding model; 399,392 peptides for presentation model		
<i>Neonmhc2</i> [61]	2019	II	CNN	MHC binding	MHC ligands	> 50,000 peptides		
<i>Neopesee</i> [226]	2018	I	ML	Immunogenicity	Epitopes with T-cell response	311 neoepitopes with T-cell response		
<i>pMTnet</i> [227]	2021	I	DL	Immunogenicity	Affinity measurements & pMHC-TCR pairs	172,422 affinity measurements & 243,747 human TCR $\beta$ CDR3 sequences & 32,607 pMHC-TCR pairs		
<i>ForestMHC</i> [228]	2019	I	RF	MHC binding	9-mer peptides from MS	> 160,000 peptides		
<i>PRIME</i> [96]	2021	I	LR	MHC binding & TCR binding	Peptides with immunogenicity or not	4958 peptides	[229]	
<i>MARIA</i> [54]	2019	II	RNN	MHC binding & protease cleavage & expression	MHC ligands	8374 peptides for presentation model & 33,909 peptides for binding model & 12,150 for cleavage model		
<i>MHCSeqNet</i> [60]	2019	I	NLP, GRU	MHC binding & ligand prediction	MHC ligands	228,348 peptides		
<i>HLAthena</i> [230]	2019	I	NN	MHC binding	MHC ligands	> 185,000 peptides		
<i>NetTCR-2.0</i> [231]	2021	I/II	CNN	Immunogenicity	TCR-peptide pairs	9204 CDR3 $\beta$ -peptide pairs & 4598 CDR3 $\alpha$ - $\beta$ -peptide pairs		
<i>NetMHCpan-4.1</i> [232]	2020	I	NN	MHC binding	Affinity measurements & MHC ligands	208,093 affinity measurements & 665,492 peptides		
<i>NetMHCIIpan-4.0</i> [232]	2020	II	NN	MHC binding	Affinity measurements & MHC ligands	108,959 affinity measurements & 381,066 peptides		

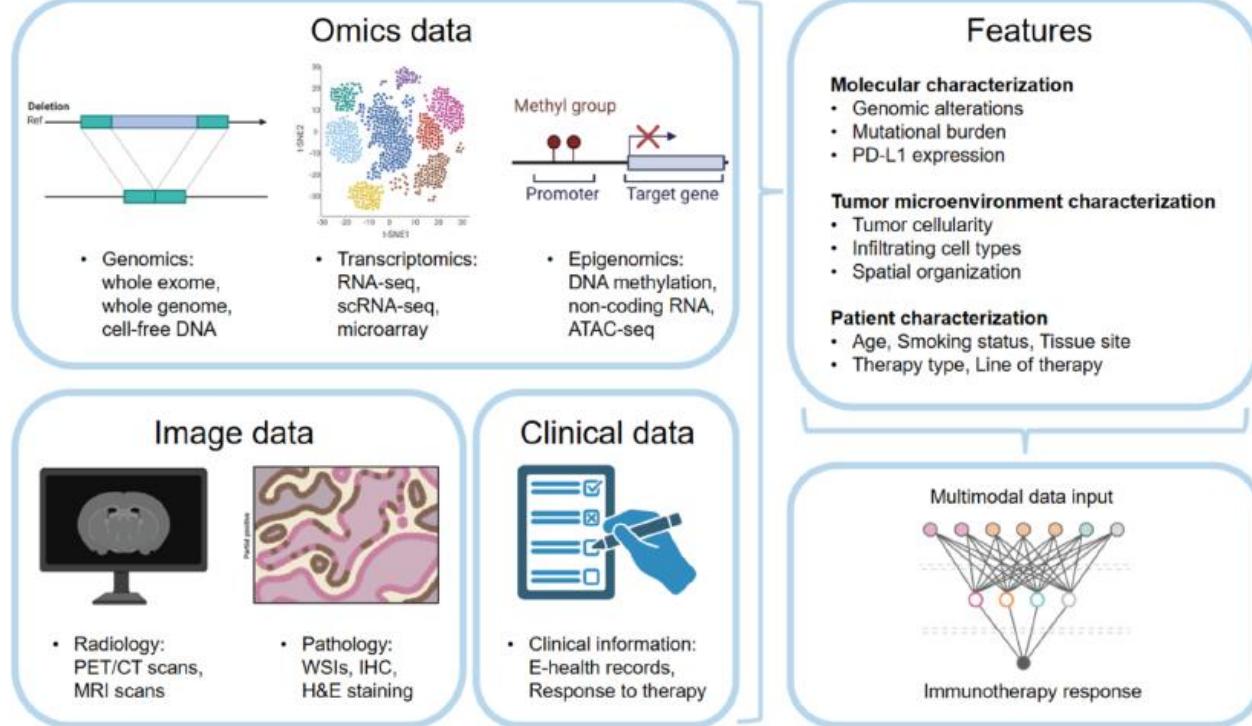


## 2. AI-based methods for antibody design

Role	Name	Model	Training input	Description	Ref.
Target binding prediction and optimization	Ens-Grad	Ensemble of CNN	51,130 CDRH3 sequences and their R2-to-R3 enrichment in phage-display panning	Designing antibody CDRH3 regions with target affinities and improved specificity	[107]
	-	CNN	CDRH3 sequences of 11,300 binding and 27,539 non-binding trastuzumab variants	Optimization of trastuzumab in affinity and pharmaceutical properties	[108]
	-	CNN and generative adversarial network	6003 non-binder and 1345 binder sequences for CTLA-4; 6052 non-binder and 1719 binder sequences for PD-1	Classifying and generating sequences of CTLA-4 and PD-1 binding antibodies	[109]
	-	LSTM	959 VH sequences	Generating antibody sequences with improved affinity to kynurenone	[110]
	-	linear discriminant analysis model	4000 sequences (2000/2000 with high and low specificity; 1516/2404 with high and low antigen binding)	Co-optimizing the affinity and specificity of emibetuzumab	[111]
	Antibody structure prediction	DeepH3	Series of CNN	1388 structures from SAbDab	Predicting inter-residue distances and orientations of antibody CDRH3 region
Antibody structure prediction	DeepAb	ReNet and bi-LSTM encoder	118,386 paired sequences from the OAS database and 1692 structures from SAbDab	Predicting relative distances and orientations of antibody Fv region and realizing structure by Rosetta	[121]
	DeepSCAb	ResNet for the inter-residue module and a transformer encoder model for the rotamer module	1433 structures from SAbDab	Predicting full Fv structures including side-chain geometries	[126]
	ABlooper	five E(n)-equivariant graph neural networks	3438 structures from SAbDab	Predicting structures of antibody CDR loops	[127]
	NanoNet	Two 2D CNN	– 2000 heavy chains of mAbs and nanobody structures	Modeling nanobodies and VH domains of antibodies	[130]
	AbLSTM	Bi-LSTM network	25,000 antibody sequences from BCR repertoire sequencing	Evaluating the nativeness of antibody candidates	[135]
	BioPhi	RoBERTa transformer encoder model for Sapiens	Human antibody repertoires in the OAS database	A platform composed of OASis which is for antibody humanness evaluation and Sapiens for antibody humanization	[136]
Pharmaceutical property	solPredict		Protein solubility data of 220 antibodies		[138]

(continued on next page)

Role	Name	Model	Training input	Description	Ref.
Pharmaceutical property	DeepSCM	Support vector machine and random forest models with language-based transfer learning 1D CNN	6596 antibody Fv sequences and their SCM scores	Predicting the apparent solubility of antibodies in histidine (pH 6.0) buffer A surrogate model for the SCM to predict antibody viscosity	[140]
	-	Random forest classifier	64 clinical-stage antibodies and their PK data	Identifying biophysical assays and in silico properties correlated with antibody PK	[144]



3. AI-based methods for immunotherapy response prediction

	Target	Summary	Data type	Tumor type	Technology	Validation	Ref.
Prediction of associated biomarker	MSI	Predict MSI directly from H&E histology	H&E histology	Gastrointestinal cancer	CNN with deep residual learning (resnet18)	AUC = 0.84	[158]
	MSI or dMMR	Identified dMMR or MSI from H&E-stained slides	H&E-stained slides	Colorectal cancer	Modified ShuffleNet/ MobileNetV2 architecture	AUPRC = 0.92–0.931, 66.6–67 % SPE and 76.0–95 % SEN	[161, 163]
	PD-L1 score	Quantitatively analyze the PD-L1 score of tumor cells	PD-L1-stained digital images	Cutaneous melanoma	ML algorithm (random forest classifier)	Correlation between label and prediction ( $r = 0.97$ , $P < 0.0001$ ).	[164]
	PD-L1 expression	Predict the TPS of PD-L1 expression	WSIs of the 22c3 assay	Non-small cell lung cancer	U-Net structure with residual blocks	ACC = 0.9326 and SPE = 0.9641	[166]
	TMB	Predicting TMB from histopathological images	H&E-stained histopathological slides	Lung adenocarcinoma	Inception-v3 and the random forest architecture	AUPRC = 0.92, precision = 0.89	[168]
	TMB status	Predicting TMB from images and clinical information	Histopathological images and clinical data	Colorectal cancer	Multi-modal deep learning model based on ResNet	AUPRC = 0.817	[172]
Deciphering TME	TME reconstruction	Prediction of 51 unique cell subpopulations	Bulk RNA-seq	Pan-cancer from TCGA	Decision tree ML deconvolution algorithm	Cytometric, immunohistochemical, or scRNA-seq	[182]
	TILs maps	Predict spatial patterns of TILs	H&E images	13 TCGA tumor types	Semi-supervised CNN and unsupervised CAE	Prediction correlate with pathologist and molecular estimates	[189]
	Tumor Cellularity	Predict tumor purity from H&E Image	Whole Slide H&E Image	Breast cancer	Weakly-supervised segmentation model with Resnet-34 CNN	Cohen's kappa coefficient = 0.69	[193]
Directly prediction	TME	Predict spatial mapping of multiple cell types	H&E images and spatial transcriptomic data	Lung adenocarcinoma		Correlated with bulk RNA-seq	[202]
	Immunotherapy response	Predict ICI treatment responses in three different cancer types	Clinical outcome and transcriptomic data	Melanoma, gastric cancer, and bladder cancer	Protein–protein interaction network (PPI)-based ML	AUCs > 0.7	[207]
	Immunotherapy response	A predictive model of immunotherapy using genomics data	Genomic data and RNA-seq	29 cancer types from TCGA	MultiModal Network	Compared to benchmark markers ( $p = 0.009$ )	[208]
	ICB efficacy	Predict ICB response	Genomic, molecular, demographic and clinical data	16 different cancer types	Ensemble learning random forest ML model	Pan-cancer AUC = 0.85	[209]
	Immunotherapy response	Prediction of response to PD-(L)1 blockade	Medical imaging, histopathological, and genomic features	Non-small cell lung cancer	Multiple-instance LR and multimodal dynamic attention	AUC = 0.80	[216]

**Table 2**

Datasets for different applications of immunotherapy.

1. Datasets for neoantigen prediction				
Name	Year	Last updated	Overview of the data	
<i>IEDB</i> [99]	2004	2023/1/30	> 1,500,000 epitopes from > 23,000 studies, including B cell, T cell, and MHC ligand assays	
<i>dbPepNeo2</i> [233]	2022	No record	801 HC neoantigens and 864,084 low LC peptidomes, validated neoantigen peptides, TCRs, and HLA peptidomes	
2. Datasets for antibody design				
Database name	Description		Link	Ref.
Observed Antibody Space (OAS)	Collecting immune repertoires		<a href="https://opig.stats.ox.ac.uk/webapps/oas/">https://opig.stats.ox.ac.uk/webapps/oas/</a>	[122]
Structural Antibody Database (SAbDab)	Database of antibody structures		<a href="https://opig.stats.ox.ac.uk/webapps/newsabdab/sabdab/">https://opig.stats.ox.ac.uk/webapps/newsabdab/sabdab/</a>	[117]
3. Datasets for immunotherapy response prediction				
Dataset name	Description		Link	Type(s) of cancer (s) Ref.
Predict MSI in gastrointestinal cancer	Histological images for MSI vs. MSS classification in gastrointestinal cancer		<a href="https://doi.org/10.5281/zenodo.2530789">https://doi.org/10.5281/zenodo.2530789</a> <a href="https://doi.org/10.5281/zenodo.2530335">https://doi.org/10.5281/zenodo.2530335</a> <a href="https://doi.org/10.5281/zenodo.2532612">https://doi.org/10.5281/zenodo.2532612</a>	Gastrointestinal cancer [150]
The Cancer Genome Atlas (TCGA)	TCGA is a pool of molecular data sets publicly accessible and freely available to cancer researchers		<a href="https://portal.gdc.cancer.gov/">https://portal.gdc.cancer.gov/</a>	Pan-cancer [161]
NCT-CRC-HE-100K and CRC-VAL-HE-7K datasets	This is a set of 100,000 image patches from H&E stained histological images of colorectal cancer (CRC) and normal tissue.		<a href="https://doi.org/10.5281/zenodo.1214455">https://doi.org/10.5281/zenodo.1214455</a>	Colorectal cancer [163]
MMDL-colorectal cancer	Predicting colorectal cancer TMB from histopathological images and clinical information		<a href="https://github.com/hkmgeneis/MMDL/tree/master">https://github.com/hkmgeneis/MMDL/tree/master</a>	Colorectal cancer [172]
Kassandra	Precise reconstruction of the TME using bulk RNA-seq and a ML algorithm trained on artificial transcriptomes		<a href="https://science.bostongene.com/kassandra/">https://science.bostongene.com/kassandra/</a>	Pan-cancer [182]
CIBERSORTx	Determining cell type abundance and expression from bulk tissues with digital cytometry		<a href="http://cibersortx.stanford.edu/">http://cibersortx.stanford.edu/</a>	Pan-cancer [183]
Scaden	DL-based cell composition analysis from tissue expression profiles		<a href="https://www.science.org/doi/10.1126/sciadv.aba2619#sec-4">https://www.science.org/doi/10.1126/sciadv.aba2619#sec-4</a>	Pan-cancer [184]
MethylCIBERSORT	Pan-cancer deconvolution of tumor composition using DNA methylation		<a href="https://www.nature.com/articles/s41467-018-05570-1#Sec12">https://www.nature.com/articles/s41467-018-05570-1#Sec12</a>	Pan-cancer [185]
MethylNet	An automated and modular deep learning approach for DNA methylation analysis		<a href="https://doi.org/10.24433/CO.6373790.v1">https://doi.org/10.24433/CO.6373790.v1</a>	Pan-cancer [186]
Human-interpretable image features (HIFs)	HIFs derived from densely mapped cancer pathology slides predict diverse molecular phenotypes		<a href="https://www.nature.com/articles/s41467-021-21896-9#data-availability">https://www.nature.com/articles/s41467-021-21896-9#data-availability</a>	Pan-cancer [194]
NetBio	Network-based ML approach to predict immunotherapy response in cancer patients		<a href="https://zenodo.org/record/4661265">https://zenodo.org/record/4661265</a> <a href="http://research-pub.gene.com/IMvigor210CoreBiologies/">http://research-pub.gene.com/IMvigor210CoreBiologies/</a> <a href="https://string-db.org/">https://string-db.org/</a>	Pan-cancer [207]
DeepTCR	DL reveals predictive sequence concepts within immune repertoires to immunotherapy		<a href="https://zenodo.org/record/6590069">https://zenodo.org/record/6590069</a>	Pan-cancer [210]
PD-(L)1 blockade in patients with NSCLC	Multimodal integration of radiology, pathology and genomics for prediction of response		<a href="https://www.synapse.org/#!Synapse:syn26642505">https://www.synapse.org/#!Synapse:syn26642505</a> <a href="https://www.ncbi.nlm.nih.gov/study/summary?id=lung_msk_mind_2020">https://www.ncbi.nlm.nih.gov/study/summary?id=lung_msk_mind_2020</a>	Non-small cell lung cancer [216]

Received: 18 October 2019

Revised: 26 November 2019

Accepted: 16 December 2019

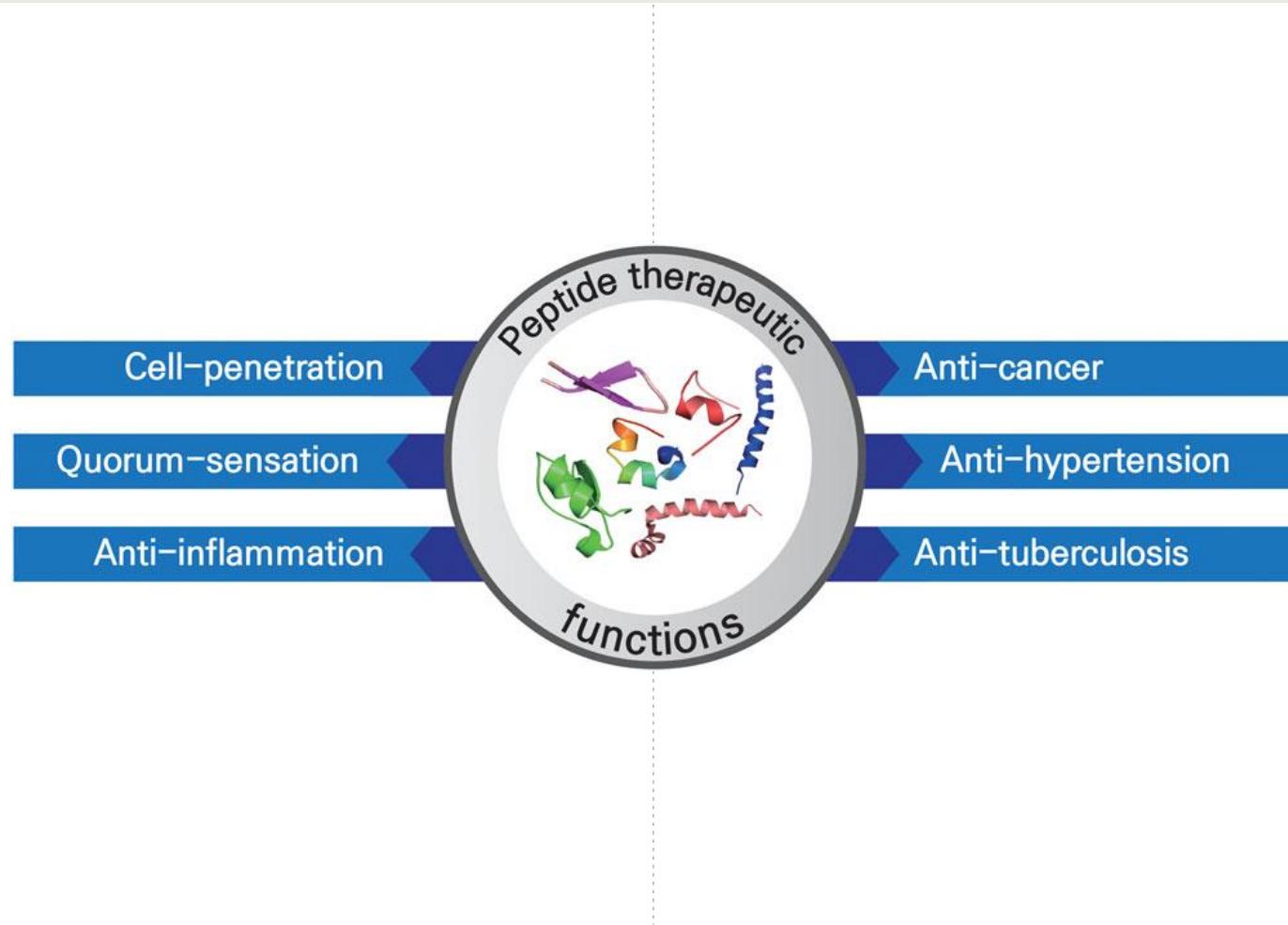
DOI: 10.1002/med.21658

**REVIEW ARTICLE**

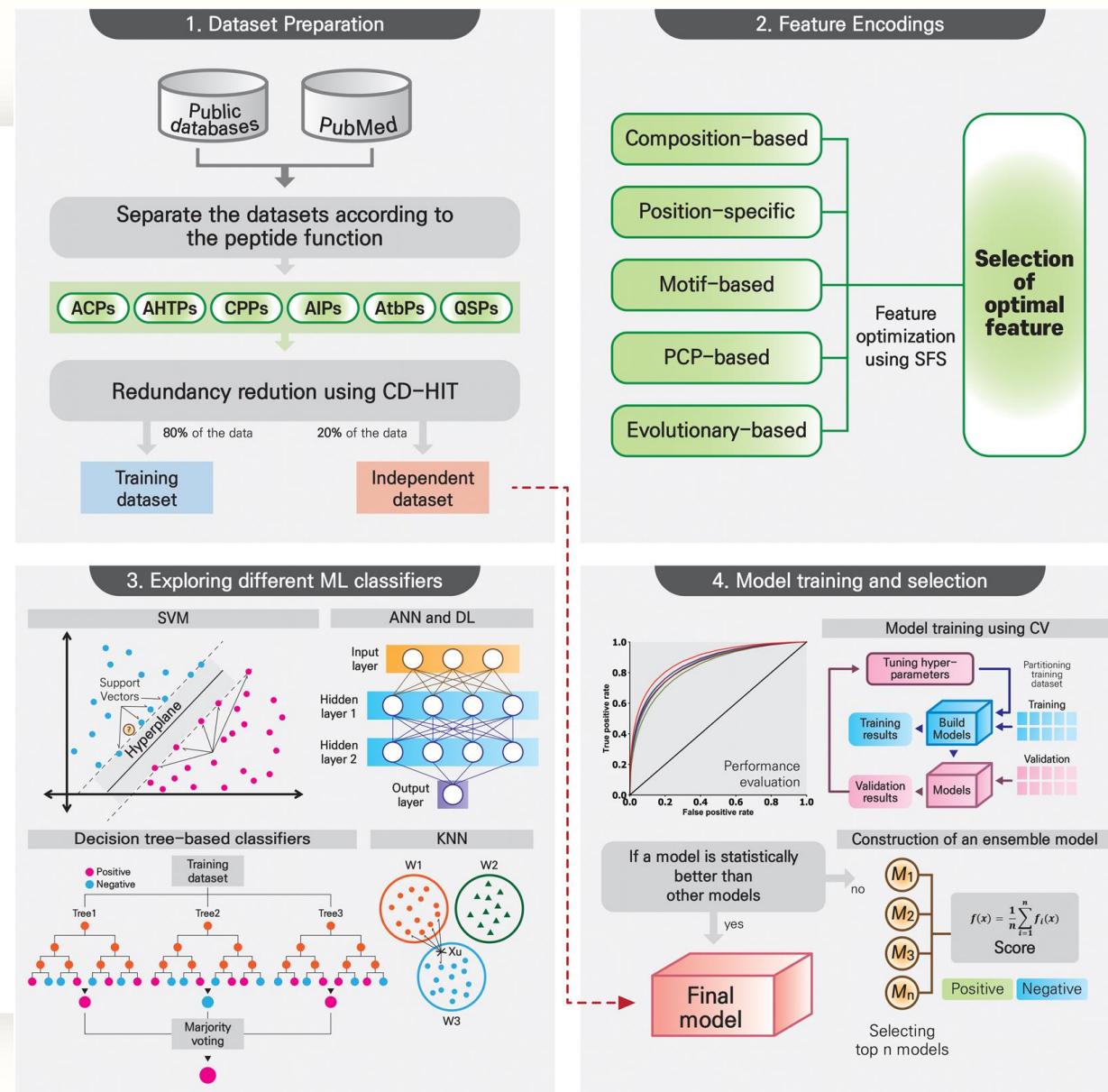
**WILEY**

# Machine intelligence in peptide therapeutics: A next-generation tool for rapid disease screening

Shaherin Basith  | Balachandran Manavalan  | Tae Hwan Shin |  
Gwang Lee



\* Source: xxxx, 2022.0x.xx



\* Source: xxxx, 2022.0x.xx

**TABLE 2** List of currently available ACP predictors evaluated in this review

Predictor	Classifier, Year	Data set size (number of positive/ negative samples)	Feature encodings	Evaluation strategy	Accuracy (CV/ independent)	Web server availability
AntiCP <sup>95</sup>	SVM, 2013	225/225 50/50	BPF	10-Fold CV	0.914/0.89	<a href="http://crdd.osdd.net/raghava/anticp/">http://crdd.osdd.net/raghava/anticp/</a>
Hajisharifi et al <sup>96</sup>	SVM, 2013	138/206 —	Local alignment kernel	5-Fold CV	0.897/—	NA
ACPP <sup>97</sup>	SVM, 2015	217/3979 40/40	Protein-relatedness measure	10-Fold CV	0.97/0.962	<sup>†</sup> <a href="http://acpp.bicpu.edu.in/predict.php">http://acpp.bicpu.edu.in/predict.php</a>
iACP <sup>98</sup>	SVM, 2016	138/206 150/150	GGDPC	5-Fold CV	0.947/0.927	<a href="http://lin-group.cn/server/iACP">http://lin-group.cn/server/iACP</a>
Li and Wang <sup>99</sup>	SVM, 2016	138/206 150/150	AAC, average chemical shifts, and RAAC	LOOCV	0.936/0.893	NA
Khan et al <sup>100</sup>	SVM, 2017	138/206 —	SAAC	LOOCV	0.933/—	NA
iACP-GAEnsC <sup>101</sup>	(SVM, RF, PNN, KNN, and GRNN), 2017	138/206 —	Amphiphilic PseAAC, GGDPC, and RAAC	LOOCV	0.965/—	NA
MLACP <sup>12</sup>	RF, 2017	187/398 422/422	AAC, ATC, DPC, and PCP	10-Fold CV	0.872/0.827	<a href="http://www.thegleelab.org/MLACP.html">www.thegleelab.org/MLACP.html</a>
SAP <sup>102</sup>	SVM, 2018	138/206 —	GGDPC	5-Fold CV	0.919/—	NA
ACPred-FL <sup>67</sup>	SVM, 2018	250/250 82/82	5 Class features	10-Fold CV	0.914/0.884	<a href="http://server.malab.cn/ACPred-FL">http://server.malab.cn/ACPred-FL</a>
mACPpred <sup>28</sup>	SVM, 2019	266/266 157/157	AAC, DPC, CTD, QSO, AAI, BPF, and CTD	10-Fold CV	0.917/0.914	<a href="http://thegleelab.org/mACPpred">http://thegleelab.org/mACPpred</a>
ACPred <sup>103</sup>	SVM, 2019	138/205 —	AAC and Am-PseAAC	LOOCV	0.956/—	<sup>‡</sup> <a href="https://codes.bio/acpred/">https://codes.bio/acpred/</a>
ACP-DL <sup>104</sup>	DL (LSTM), 2019	376/364 —	BPF, k-mer sparse matrix	5-Fold CV	0.815/—	<sup>§</sup> <a href="https://github.com/haichengyi/ACP-DL">https://github.com/haichengyi/ACP-DL</a>

**TABLE 7** List of currently available CPP predictors evaluated in this review

Predictor	Classifier, Year	Data set size (number of positive/ negative samples)	Feature encodings	Evaluation strategy	Accuracy (CV/ independent)	Web server availability
Dobchev et al <sup>61</sup>	ANN, 2010	49/10 23/2	QSAR descriptors	3-Fold CV	0.983/0.96	NA
Sanders et al <sup>62</sup>	SVM, 2011	111/34 —	PCP	10-Fold CV	0.917/—	NA
CellPPD <sup>63</sup>	SVM, 2013	708/708 99/99	A combination of motif and BPF	5-Fold CV	0.974/0.813	<a href="http://crdd.osdd.net/raghava/cellppd/">http://crdd.osdd.net/raghava/cellppd/</a>
CPPpred <sup>66</sup>	N-to-1 NN, 2013	74/100 47/47	Motif	5-Fold CV	0.776/0.830	<sup>†</sup> <a href="http://bioware.ucd.ie/cpppred">http://bioware.ucd.ie/cpppred</a>
C2Pred <sup>115</sup>	SVM, 2016	411/411 —	DPC	5-Fold CV	0.836/—	<a href="http://lin-group.cn/server/C2Pred">http://lin-group.cn/server/C2Pred</a>
SkipCPP-Pred <sup>67</sup>	RF, 2017	462/462 —	Adaptive k-skip-2-gram	LOOCV	0.906/—	<a href="http://server.malab.cn/SkipCPP-Pred/Index.html">http://server.malab.cn/SkipCPP-Pred/Index.html</a>
CPPred-RF <sup>91</sup>	RF, 2017	462/462 —	PseAAC, DPC and PCP	LOOCV	0.916/—	<a href="http://server.malab.cn/CPPred-RF">http://server.malab.cn/CPPred-RF</a>
MLCPP <sup>93</sup>	ERT, 2018	427/427 311/311	AAC and PCP	10-Fold CV	0.883/0.896	<a href="http://www.thegleelab.org/MLCPP">www.thegleelab.org/MLCPP</a>
KELM-CPPpred <sup>116</sup>	KELM, 2018	408/408 96/96	AAC and the motif-based features	10-Fold CV	0.870/0.870	<a href="http://sairam.people.iitgn.ac.in/KELM-CPPpred.html">http://sairam.people.iitgn.ac.in/KELM-CPPpred.html</a>
CPPred-FL <sup>92</sup>	RF, 2018	462/462 —	19 Probabilistic features	10-Fold CV	0.921/—	<a href="http://server.malab.cn/CPPred-FL">http://server.malab.cn/CPPred-FL</a>
G-DipC <sup>117</sup>	XGB, 2019	1223/1223 —	General DPC	5-Fold CV	0.985/—	NA
PEPred-SUITE <sup>16</sup>	RF, 2019	370/370 92/92	15 Probabilistic features	10-Fold CV	0.912/0.880	<a href="http://server.malab.cn/PEPred-Suite">http://server.malab.cn/PEPred-Suite</a>

## Deep Learning Enables Discovery of a Short Nuclear Targeting Peptide for Efficient Delivery of Antisense Oligomers

Eva M. López-Vidal,<sup>▼</sup> Carly K. Schissel,<sup>▼</sup> Somesh Mohapatra, Kamela Bellovoda, Chia-Ling Wu, Jenna A. Wood, Annika B. Malmberg, Andrei Loas, Rafael Gómez-Bombarelli,\* and Bradley L. Pentelute\*



Cite This: *JACS Au* 2021, 1, 2009–2020



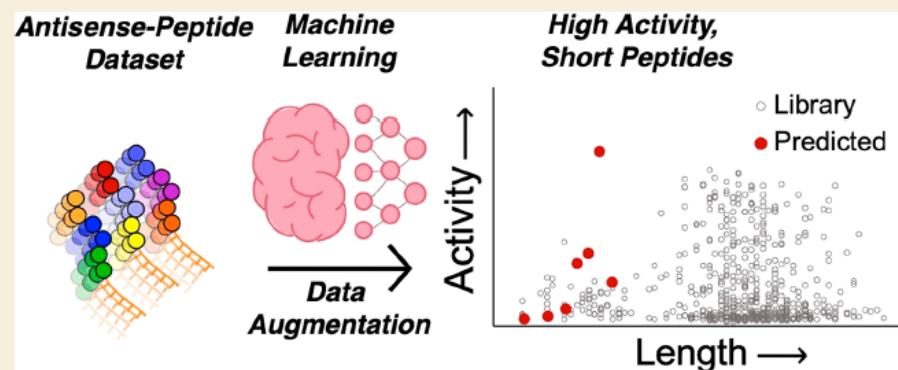
Read Online

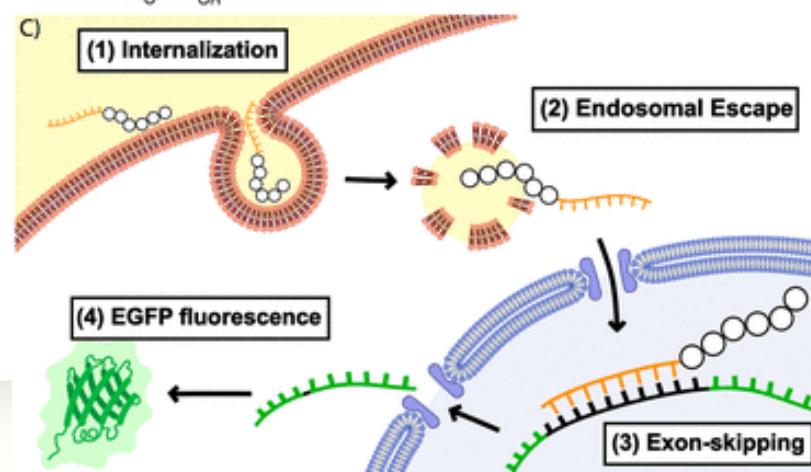
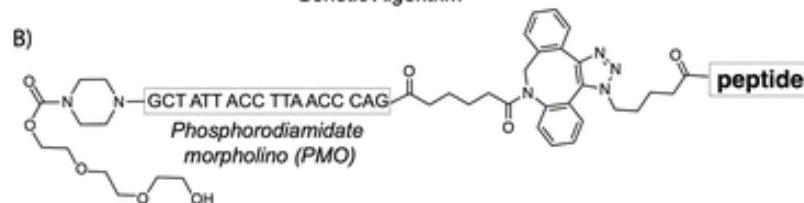
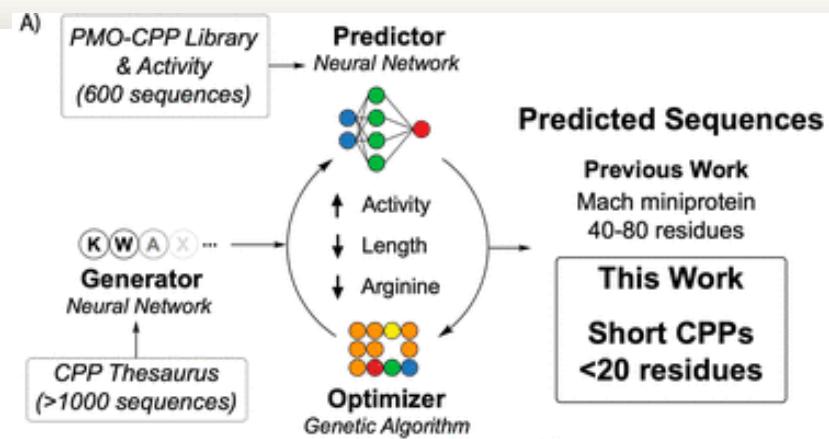
ACCESS |

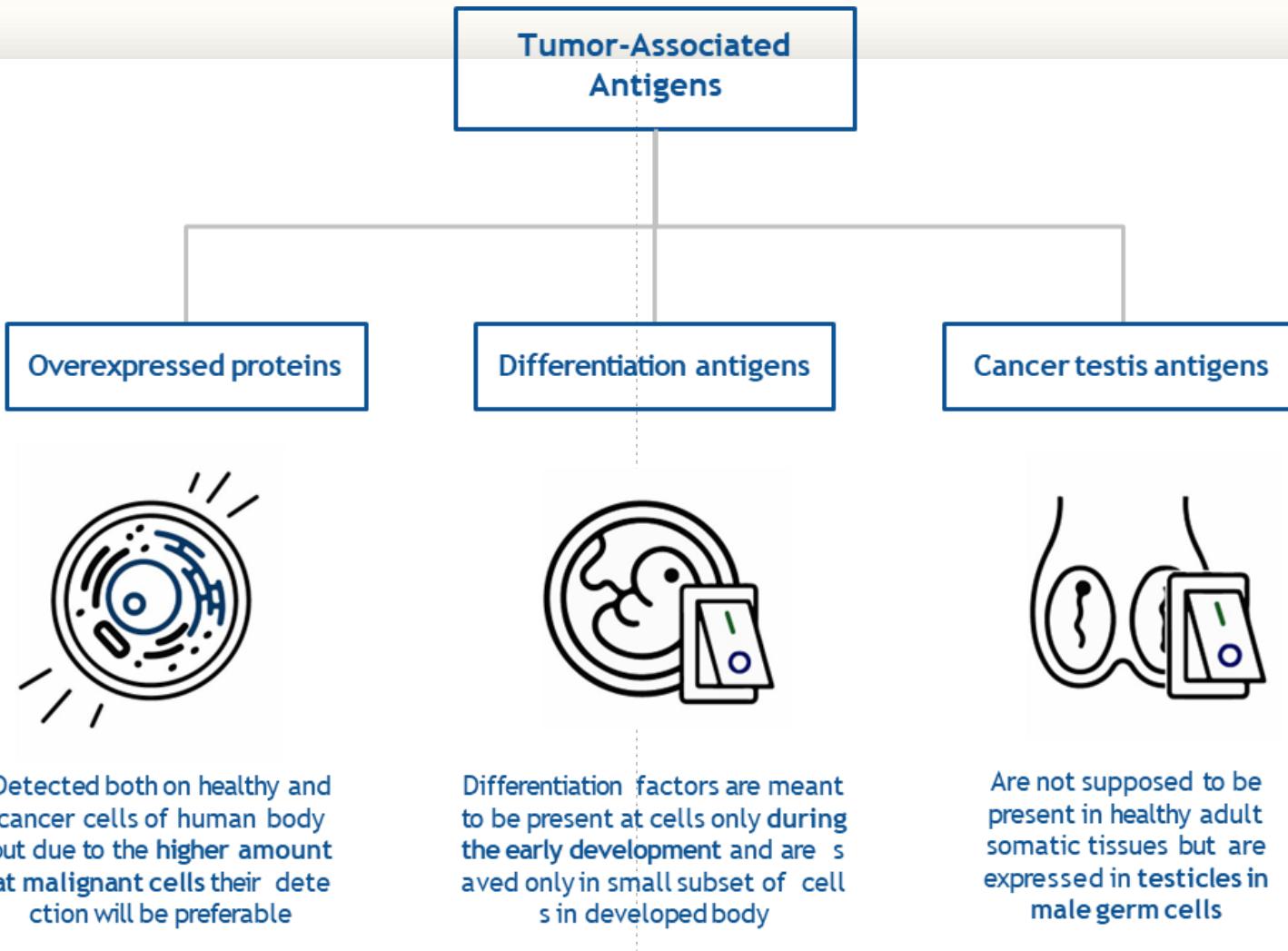
Metrics & More

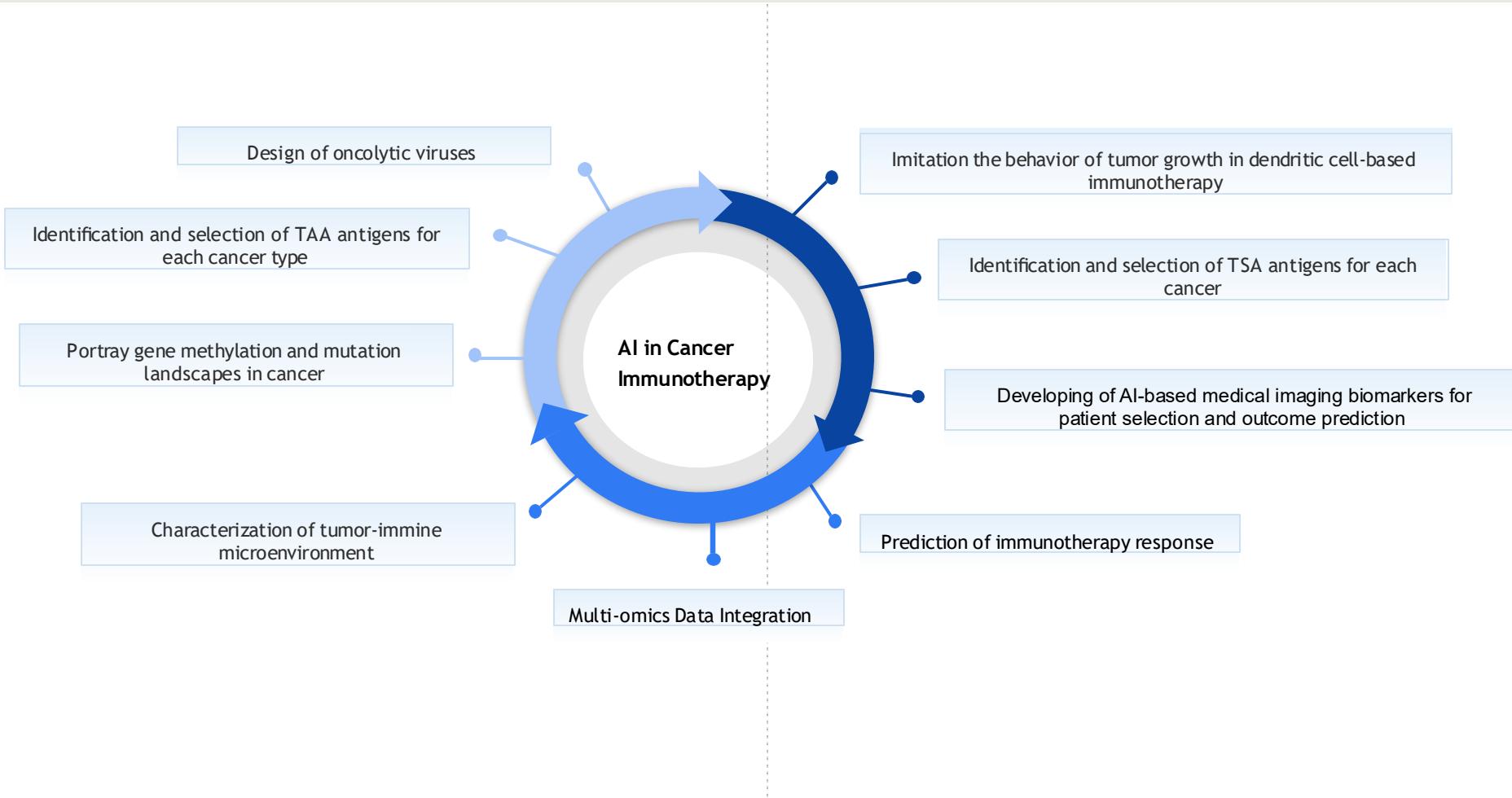
Article Recommendations

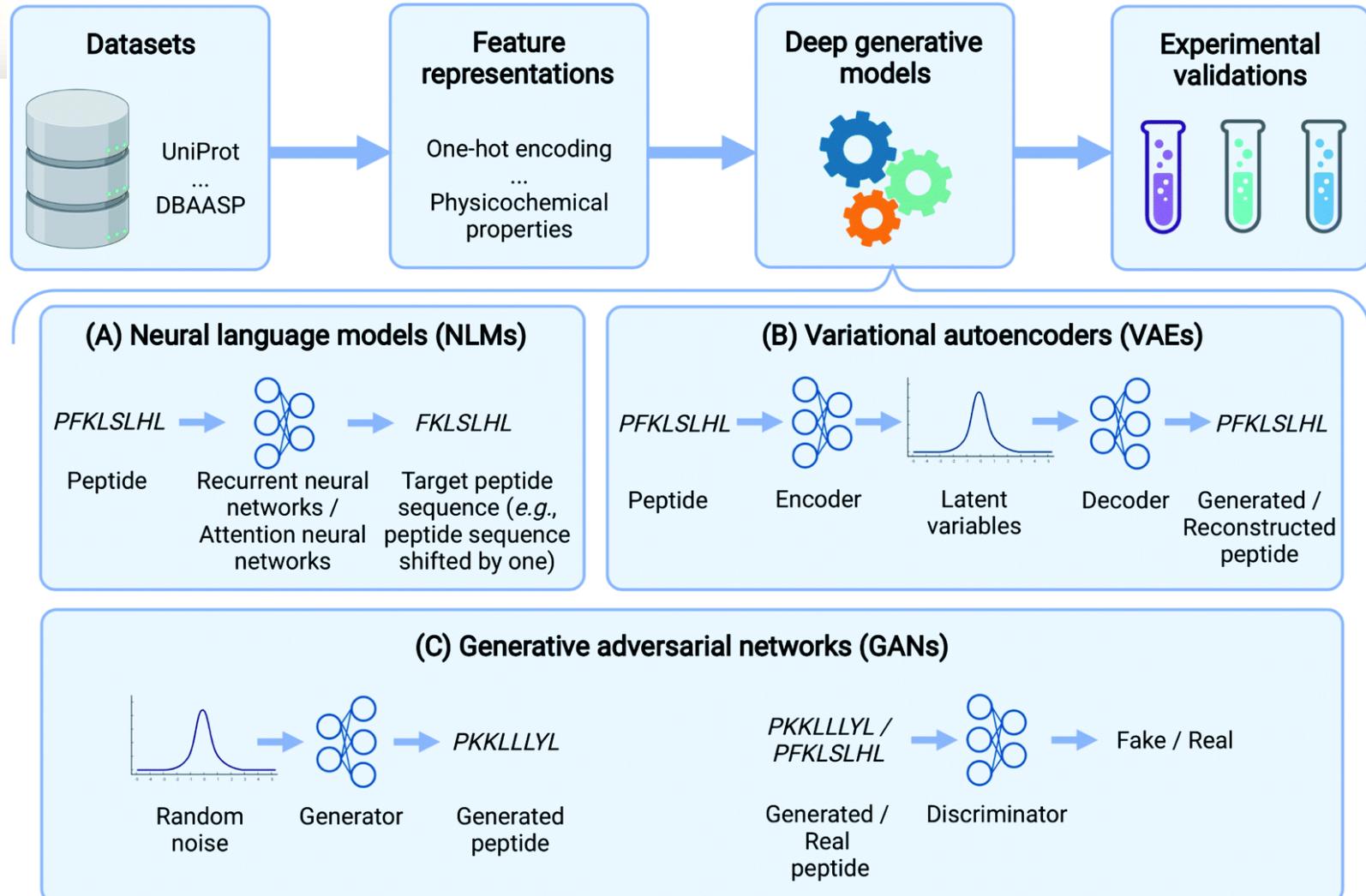
Supporting Information











# 실습

- 실습 코드: [https://github.com/fourmodern/2025\\_aidrugdiscovery](https://github.com/fourmodern/2025_aidrugdiscovery)
- 구글 코랩: <https://colab.research.google.com/>
- 구글 AI Studio: [https://aistudio.google.com/prompts/new\\_chat](https://aistudio.google.com/prompts/new_chat)