

Natural Language Processing Methods for the Study of Protein-Ligand Interactions

James Michels,[†] Ramya Bandarupalli,[‡] Amin Ahangar Akbari,[‡] Thai Le,[¶]

Hong Xiao,^{*,†} Jing Li,^{*,‡} and Erik F. Y. Hom^{*,§}

[†]*Department of Computer Science, University of Mississippi, University, MS*

[‡]*Department of BioMolecular Sciences, School of Pharmacy, University of Mississippi, University, MS*

[¶]*Department of Computer Science, Indiana University, Bloomington, IN*

[§]*Department of Biology and Center for Biodiversity and Conservation Research, University of Mississippi, University, MS*

E-mail: hxiao1@olemiss.edu; jli15@olemiss.edu; erik@fyhom.com

Abstract

Natural Language Processing (NLP) has revolutionized the way computers are used to study and interact with human languages and is increasingly influential in the study of protein and ligand binding, which is critical for drug discovery and development. This review examines how NLP techniques have been adapted to decode the “language” of proteins and small molecule ligands to predict protein-ligand interactions (PLIs). We discuss how methods such as long short-term memory (LSTM) networks, transformers, and attention mechanisms can leverage different protein and ligand data types to identify potential interaction patterns. Significant challenges are highlighted, including the scarcity of high-quality negative data, difficulties in interpreting model decisions, and sampling biases of existing datasets. We argue that focusing on improving data quality, enhancing model robustness, and fostering both collaboration and

competition could catalyze future advances in machine-learning-based predictions of PLIs.

1. Introduction

Proteins play a pivotal role as molecular machines essential for biological function. Their functionality often depends on site-specific binding interactions with small molecule ligands that cannot be studied within the protein-protein interaction framework.¹ Understanding such protein-ligand interactions (PLIs) is central to drug discovery and development^{2,3} as well as protein engineering efforts.^{4,5} Although laboratory experimentation is the traditional approach to studying PLIs and generating "ground-truth" data for a specific protein-ligand pair, it is both costly and time-consuming, often taking weeks to months.⁶ Computational approaches that simulate the underlying physics and chemistry of PLIs such as molecular docking⁷ or dynamics simulations⁸ can be less resource intensive but nevertheless demand significant computational and time investment.⁹

In recent years, machine learning (ML) has provided new avenues for analyzing biological data, leveraging statistical and algorithmic techniques to distill potentially human-interpretable insights with little manual intervention. ML models have successfully predicted various protein and molecular attributes,¹⁰⁻¹⁴ and have moved us closer to solving the protein folding problem.^{15,16} As the excitement for ML use in the biological sciences grows, the prediction of protein-ligand interactions appears increasingly possible given recent advances in both ML and Natural Language Processing (NLP),^{17,18} the computational study of language.¹⁹

1.1. Overview of Natural Language Processing (NLP)

NLP centers on the computational analysis and manipulation of language constructs to bridge the gap between human communication and computer automation. NLP has ex-

perienced significant recent breakthroughs as demonstrated by the proliferation of widely used chatbots such as OpenAI’s ChatGPT,^{20,21} Anthropic’s Claude,²² and Microsoft’s Bing Copilot.²³ NLP has been further used to summarize texts, deduce author sentiment, solve symbolic math problems, and even generate programming code.^{24–27} The effectiveness of NLP is predicated on (human) languages having a structured symbolic syntax and set of rules to assemble basic units known as “tokens” (e.g., characters, words, or punctuation) to form higher-order constructs such as sentences or paragraphs. The structured outputs of such a system reflect the grammar, conventions, and styles of the associated language. In NLP, tokens are transformed to encode “meanings” through mathematical vectors such that tokens of similar meaning are positioned closer together in the representational vector space. By analyzing a large collection of data, NLP methods aim to infer emergent relationships between tokens that define the “rules” of a language. Importantly, this inferred set of rules can then be used to perform predictive tasks such as separating tokens into categories, translating text from one language to another, and even predicting whether a literary work will see commercial success.^{28–30}

In the biological domain, NLP methods have been used for a variety of predictive tasks, including inferring disease-gene associations,³¹ predicting tumor gene expression patterns,³² and assigning functional annotations to various protein-coding genes.¹² More recently, NLP has been applied with unprecedented success in DeepMind’s AlphaFold algorithm to predict three-dimensional protein structures given only protein sequence data.^{15,16} Despite impressive advances, the creation of these NLP models is associated with a sizable computational burden (see for example,^{33–38}) and it remains a challenge to understand what and which specific features of the input sequence data fuel predictive success.

Below, we review contemporary NLP methods as they have been applied in the study of PLIs in recent years. We first describe the relationship between common protein and ligand text representations vis-à-vis the characteristics of human language. We then discuss the dominant NLP-based approaches used to study PLIs and provide a comprehensive overview

of the diversity of ML models used in this space. We conclude with reflections on remaining challenges in the field and areas that merit future development.

2. The "Language" of Proteins

Protein sequences are akin to human language in that they possess a hierarchical order of construction and embody embedded information (Fig. 1). Human language text is inherently ordered with characters of an alphabet assembled linearly and grouped into words, phrases, and sentences that convey an emergent message. Protein sequences similarly obey a hierarchy of assembly, with amino acids (AAs) serving as the alphabet. When AAs are strung together, secondary structural motifs, domains, and quaternary (multi-domain-interacting) structures may emerge with properties that contribute to function.^{39,40} While external factors such as post-translational modifications and cellular state can play a substantial role in dictating protein three-dimensional structure and function, the AA sequence represents the essential blueprint that ontologically defines the properties of a protein.⁴¹⁻⁴³ This fact has served as the foundation for bioinformatic analysis of proteins.⁴⁴ Individual AAs and common subsequences contribute to the "information" of the overall protein just as words contribute to the meaning of a text.

3. The "Language" of Ligands

The chemical structures of molecules can be similarly translated into text-based notations and analyzed computationally.⁴⁵ However, unlike the elements of human text and protein sequences, the chemical connectivity patterns of small molecule ligands are not one-dimensional. Nevertheless, text-based schema has been developed to represent chemical information in a manner convenient for computational analysis,⁴⁶ with the Simplified Molecular-Input Line-Entry System (SMILES) format being one of the most widely used.⁴⁷

SMILES strings are text representations constructed over a depth-first traversal of a

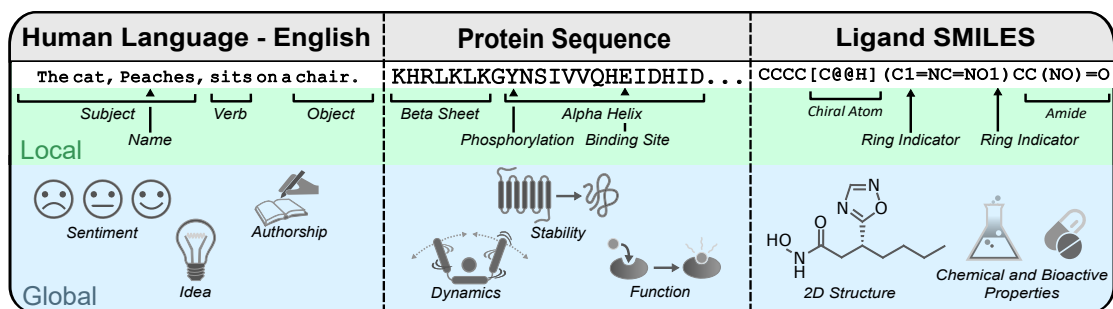


Figure 1: The Language of Protein Sequences and SMILES: NLP methods can be applied to text representations to infer local and global properties of human language, proteins, and molecules alike. Local properties are inferred characteristics of sub-sequences in text: (i) for a human language, this can include part of speech or a role a specific word serves; (ii) for a protein sequence, this can include secondary structures, post-translational modifications, and functional sites; (iii) for a SMILES string, this can include functional groups and characters used within SMILES syntax to indicate chemical attributes. Global properties are inferred from a text in its entirety: (i) for a human language, this can include information such as authorship, tone, and synopses; (ii) for a protein sequence, this can include the protein’s structure, stability, and dynamic properties; and (iii) for a SMILES string, this can include the ligand’s 2D molecular structure and other biochemical properties.

two-dimensional molecular graph (Fig. 1), with atoms, atomic properties, bonds, and structural properties represented by characters following an established set of conversion rules. Given the memory-efficient and somewhat human-readable format of SMILES, it has become a standard in chemical databases and computational tools^{48–50} and the most commonly used text representation in PLI studies. Although SMILES lacks an intuitive way to determine a chemical equivalent of a “word”, there is a well-defined grammar to denote properties and substructures of a molecule. Moreover, the same molecule can be represented by multiple different SMILES strings,⁴⁷ which is similar to how there could be multiple sentence constructions to convey the same idea in human languages. In NLP applications, incorporating tokens with the same meaning into the training process can yield a robust predictive model.⁵¹ The use of multiple SMILES per molecule has been leveraged to guide ML models to discern which parts of a ligand contribute to drug potency.⁵²

4. Protein-Ligand Interaction Data and Datasets

Protein-ligand binding is a complex process dictated by many factors including protein states, hydrophobicity/hydrophilicity, and conformational flexibility.⁵³ The question of *how* to represent a protein and ligand in a computational space is critical and multifaceted. A wealth of information has been collected experimentally and generated through simulation studies on the properties of proteins and ligands, but these data are highly variable with regard to type, quality, and quantity. This section catalogs several primary data representations used in PLI studies. We also discuss the availability, selection, and curation of available data for machine-learning-based training and evaluation.

Protein and ligand representations are typically sequence- or structure-based. Unlike sequence-based text formats, structure-based information can appear in multiple forms, e.g., atomic coordinates of protein-ligand complexes or contact maps. Some structural information can be artificially reconstructed from sequence-based formats through algorithms such as AlphaFold for proteins¹⁵ and RDKit for ligands.⁵⁴ PLI studies using machine-learning methods will typically select either sequence-based or structure-based inputs, although there is a growing use of mixed input data types.^{55,56} For example, a mixed-data study may represent proteins by AA sequences but ligands by atomic coordinates, a choice based in part on the fact that highly accurate 3D chemical structures are easier to obtain than those of proteins and that full-atom representations of ligands are not memory intensive.

Other data can also be incorporated to augment ground-truth information about PLIs. For example, molecular weights, polarity, and bioactive properties can be incorporated into models to further improve the prediction of PLIs.^{57,58} Studies that included molecular weights, ligand polar surface area, and protein aromaticity,⁵⁷ or bioactive properties of chemical and clinical relevance⁵⁸ have resulted in improved predictions of binding affinity. Leveraging multiple-sequence alignment or phylogenetic information to identify co-evolutionary trends among AAs and sites of covalent modification has been shown to dramatically improve the accuracy of structural predictions of protein-ligand complexes.¹⁶ The use of

non-sequence/non-structural data can enable models to yield better predictive performance for characterizing protein and ligand and their interactions than models that do not.⁵⁷

Given a protein-ligand representation, several predictive tasks are possible. *Classification* studies seek to categorize PLIs into distinct groups, for example, whether a protein-ligand pair binds or not. These models are relatively simple and allow for input from various sources. *Regression* studies use a continuous functional metric to characterize PLIs such as a binding affinity/dissociation constant (K_d) or inhibition constant ($IC50$). Continuous target variables allow for the involvement of numerical values derived directly from 'ground-truth' experimental data in both training and evaluation. Databases like PDBBind⁵⁹ contain functional metrics such as K_d , and $IC50$, but not all protein and ligand pairings cataloged have such metrics available, for example, complexes identified from X-ray crystallography, Cryo-EM, or NMR screening studies.^{6,60} Since regression studies require quantitative PLI data and not merely whether a protein and ligand interact, relevant dataset sizes may be smaller than those for classification. However, gathering such data is a laborious process in terms of both time and laboratory resources.

Data for the study of PLIs can be manually curated by domain experts or sourced from existing datasets, such as PDBBind⁵⁹ and the Directory of Useful Decoys-Enhanced (DUD-E),^{61,62} which includes tens of thousands of diverse pairings. Other datasets such as the Davis⁶³ and KIBA⁶⁴ datasets of kinase inhibitors, focus on particular types of proteins. While pre-assembled datasets are tempting to use out of convenience, relevant data need to be selected with an intended predictive task in mind. Table 1 contains a collection of existing PLI datasets and databases for consideration.

5. Machine Learning and NLP for PLIs

The general workflow for any ML-based study can be broadly characterized into three stages: data preparation, model creation, and model evaluation. A visual aid summarizing these

Table 1: Datasets and Databases for PLI Prediction

Dataset Name	Year	Proteins	Ligands	Interactions	Protein Category	Ligand Category
<i>Functional Data Available</i>						
Protein Data Bank (PDB) ⁶⁵	2000	220,777	-	-	General	General
brenda ⁶⁶	2002	8,423	38,623	-	Enzyme	General
Natural Ligand Database (NLDB) ⁶⁷	2016	3,248	-	189,642	General	General
DrugBank ⁶⁸	2006	4,944	16,568	19,441	Human Proteome	General
PDBBind ⁵⁹	2004	-	-	23,496	General	General
BindingDB ⁶⁹	2007	2294	505,009	1,059,214	General	General
Davis-DTA ⁶³	2011	442	68	30,056	Kinases	Kinase Inhibitors
ChEMBL ⁷⁰	2012	15,398	2,399,743	20,334,684	General	General
DUD-E ⁶²	2012	102	22,886	2,334,372	General	General
PSCDB ⁷¹	2011	-	-	839	Human Proteome	General
Iridium Database ⁷²	2012	-	-	233	General	General
KIBA ⁶⁴	2014	467	52,498	246,088	Kinases	Kinase Inhibitors
dbHDPLS ⁷³	2019	-	-	8,833	General	General
PDID ⁷⁴	2016	3,746	51	1,088,789	Human Proteome	General
BindingMOAD ⁷⁵	2020	11,058	20,387	41,409	General	General
CovPDB ⁷⁶	2022	733	1,501	2,294	General	General
PSnpBind ⁷⁷	2022	731	32,261	640,074	General	General
PLAS-5k ⁷⁸	2022	-	-	5,000	Enzyme	-
Protein-Ligand Binding Database (PLDB) ⁷⁹	2023	12	556	1831	Carbonic Anhydrase, Heat Shock Protein	General
BioLiP2 ⁸⁰	2023	426,209	-	823,510	General	General
<i>Functional Data Unavailable</i>						
Database of Interacting Proteins ⁸¹	2004	28,850	-	81,923	Different Species	-
Protein Small-Molecule Database ⁸²	2009	4,916	8,690	-	General	General
CavitySpace ⁸³	2022	23,391	-	23,391	General	General

Note: Published datasets may provide periodic updates in the future. Datasets marked with the “functional data available” label contain continuous metrics.

-: Exact information is either not included in the source or is not readily obtainable.

+: Protein-ligand complexes are available with the dataset.

stages is presented in Figure 2. For PLI studies, data preparation typically entails selecting the types and formats of protein and ligand data (e.g., sequence and/or structural). ML model creation may involve the following three tasks, although the boundary between these tasks could be fuzzy at times: (i) *Extract*: the “extraction” of vector “embeddings” from the protein and ligand input data, which can be used in computational operations (described in Section 5.2) (ii) *Fuse*: the fusion of protein and ligand vector embeddings, and (iii) *Predict*: the prediction of a PLI target property as a model’s output. The predictive capability of

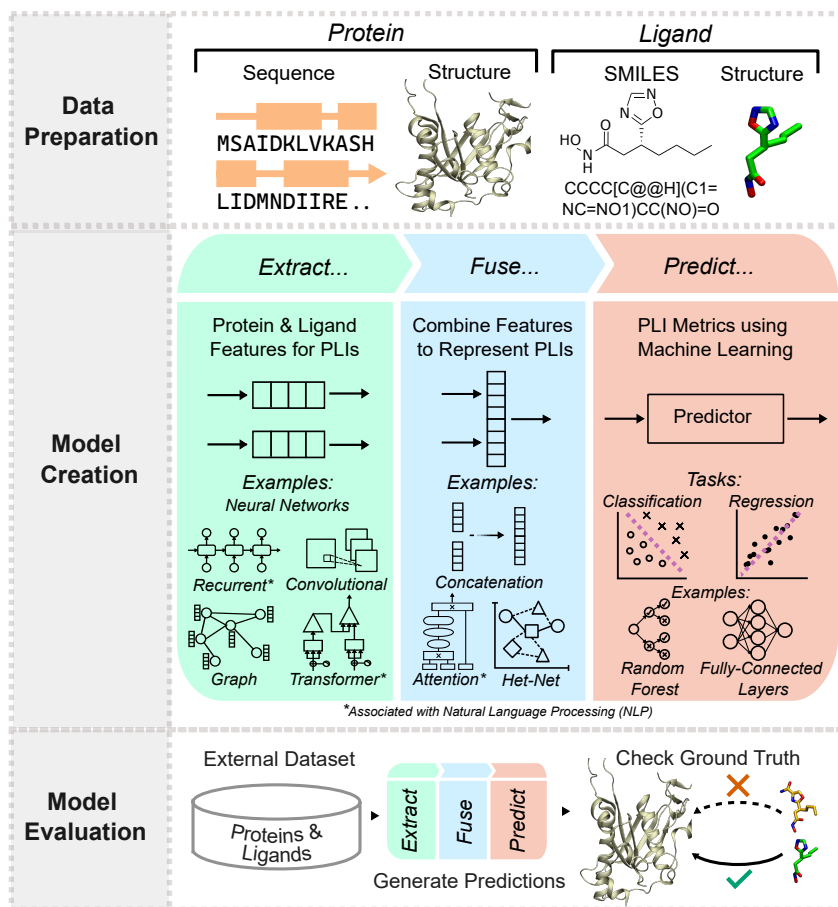


Figure 2: Summary of the Data Preparation, Model Creation, and Model Evaluation Workflow. Model Creation for PLI studies follows an Extract-Fuse-Predict Framework: input protein and ligand data are extracted and embedded, combined, and passed into a machine learning model to generate predictions.

the model would be ideally validated against results from other studies and/or real-world measurements in a model evaluation stage. While data preparation and extraction steps have typically been the focus of most research efforts, every component of the workflow is crucial to successful PLI predictions.

5.1. The Extract-Fuse-Predict Framework

A variety of models for PLI prediction have been constructed in recent years, and these models tend to fall into four general categories: (1) *sequence-based*, where protein sequences and SMILES are used to represent protein and ligand, respectively; (2) *structure-based*,

where structural information is included in the representation of both protein and ligand; (3) *mixed representations*, where both structural and sequence information are involved; and (4) *sequence-structure-plus*, which substantially incorporates other ground-truth information beyond sequence and structural data (such as molecular weights or polar surface area⁵⁷). Tables 2 through 5 summarize several representative NLP-based PLI prediction studies across these categories over the past five years. Although PLI studies could be categorized in other ways—for example by the ML model used (neural network, decision tree, etc.) or by the predictive task type (classification vs. regression)—we have chosen to emphasize a categorization based on input data type since the computational methods used for sequence text and structural data comprise a major difference.

5.2. Extraction of Embeddings

NLP approaches deconstruct text into individual tokens or "units of meaning" for use in computational operations and inferences via a process referred to as "tokenization". Schema for tokenization, aside from character-based and word-based, can also be sub-word-based. Sub-word-based tokenization breaks down text into units smaller than words to create a wider vocabulary; it is commonly selected when the definition of a "word" is unclear, as sub-words can be used as a means to discover "words".^{84,85} Common ways to assemble sub-words include methods such as "n-grams", where each sub-word has a select fixed-length value n (e.g., "Sma", "mar", "art", etc. for $n=3$). While sub-word tokenization has been attempted in PLI studies for both protein (e.g., amino-acid k-mers such as "KHR", "LKL", "KGY") and ligand (e.g., "CCCC", "[C@@H]"),⁸⁶⁻⁹⁰ the current trend is to use amino acids and/or individual atoms directly as tokens.

To be processed computationally, tokens must be translated into a numerical form through a process known as "embedding". There are many types of token embedding, but generally speaking, they are designed to capture either a particular token meaning, frequency, or both^{91,92} and represented by a multidimensional vector. The direction of a token's vector

embedding effectively represents its “meaning” while its magnitude represents the strength by which that meaning is conveyed. In isolation, each token could possess multiple meanings (e.g., the word ”run” has multiple meanings⁹³), and so context may be necessary to impart an intended meaning. NLP has been demonstrated to be highly effective at extracting patterns that convey context-dependent meanings from a large corpus of text. Embeddings that capture semantic meaning and relationships can then be used for many other tasks aside from predicting whether a protein interacts with a ligand, such as predicting protein and ligand solubilities.^{94,95}

Token embedding is typically accomplished using a neural network (NN) architecture that approximates nonlinear relationships between the “inputs” of the network (the data) and its “outputs” (the predictions).⁹⁶ Neurons in an artificial NN receive, integrate, and transmit signals to other neurons through a nonlinear response function and are arranged in layers. Information is passed from an input layer through one or more intermediate ”hidden” layers to an output layer. Interconnection weights that govern the strength of influence of one neuron on another are crucial parameters of an NN. A wide variety of NNs have been applied to studying PLIs although not all are commonly used in NLP. Nevertheless, two types of NNs are commonly associated with NLP: Recurrent Neural Networks (RNNs) and attention-based NN models.⁹⁷ Below, we highlight the details necessary to understand how RNNs, attention, and other non-NLP-driven NNs have been used to glean global patterns essential for PLI predictive tasks.

5.2.1. Recurrent Neural Networks

RNNs are specialized for processing sequential data in which the order of the data matters. Consider an input data sequence in which individual tokens are ordered by a time-step and embody a particular yet unknown pattern over the length of the sequence. In traditional NNs, information flows from the input layer to the output in a single pass, making it difficult to decipher any interdependencies between earlier tokens and subsequent ones. To remedy

this, the RNN architecture introduces recurrence units through which the processing of the input sequence at the current time-step will also update the "hidden states" that nonlinearly capture the information of all input tokens up to the current time-step. These hidden states are functionally equivalent to the hidden layers of traditional NNs but differ by updating *recurrently*, where information is carried over from previous time-steps to the current time-step. Thus, the dependencies between tokens of the sequential inputs can be captured. For example, given a protein sequence for which each AA is a token, an RNN would process the sequence of AAs one at a time to create and maintain a mapping for the next AA in the sequence accounting for all input tokens seen so far. While effective in many NLP tasks, early RNNs commonly suffered diminishing returns with increasing text length. This was due to systematic and nondiscriminatory retention of information from *all tokens*, including outlier tokens that contribute little informationally to the underlying pattern.

To minimize diminishing returns, RNNs were modified to Long Short-Term Memory (LSTM) models. The signature component of LSTMs is the "forget gate", which selectively inhibits information not concordant with previously learned patterns found from processing prior tokens.⁹⁸ For example, in the task of predicting secondary structures, an LSTM's forget gate can attenuate the contribution of AAs that do not correlate with any defined secondary structural element.^{99,100} Bidirectional LSTMs (BiLSTMs) have also been developed to capture both preceding and subsequent tokens in a sequence string by applying an LSTM to text in both original and reverse order, and concatenating each of the resulting embeddings end-to-end.¹⁰¹

LSTMs and BiLSTMs are promising embedding approaches for predicting binding affinities of proteins and ligands.¹⁰²⁻¹⁰⁴ However, their effectiveness has been limited by the size of the dataset that the LSTM/BiLSTM architecture can efficiently process. Most successful applications of LSTM to date have been applied to only relatively small training datasets, on the order of a few thousand proteins and ligand pairs. This limitation mainly arises from the inherently non-parallel design, which makes training on large datasets slow and compu-

tationally expensive. Thus, NN architectures that leverage parallelization will be important to ensure reasonable training and prediction runtimes.

5.2.2. Attention-Based Architectures

Protein lengths can vary dramatically, from Insulin with 51-AAs to "giant proteins" that can exceed 85,000 AAs.¹⁰⁵ To use large amounts of sequence data to effectively process and predict PLIs for which long-distance interactions may be impactful, several alternatives to RNN have been proposed. The "neural attention"—or simply "attention"—mechanism is an important recent breakthrough by which "attention weights" are dynamically calculated to quantify the relative contribution of different input tokens or elements to a predictive end goal. In many NN architectures, attention can also incorporate hidden states into the calculation, allowing a more sophisticated mechanism for capturing longer-range correlations in deeper layers.

Attention mechanisms have proven highly compatible with traditional protein sequence analysis approaches in identifying long-distance interactions between AAs of a protein.¹⁰⁶ In PLI studies, attention mechanisms can dynamically adjust the contribution of specific AAs or ligand atoms to a predictive outcome by amplifying interaction sites with higher attention scores and downplaying less relevant ones. This process mirrors the biological intuition that certain residues and atoms are more critical for binding in a protein-ligand complex than others. The use of attention mechanisms has enabled the identification of AAs in proteins and atoms in a ligand that are highly cross-correlated and appear to physically interact,¹⁰⁷ although the degree of success in identifying physically interacting sites remains to be carefully assessed. Attention has also provided an effective way to "fuse" protein and ligand representations in binding prediction models.^{56,86,104,108,109} Figure 3 presents an example of how attention weights have been used to reveal potentially interacting sites between a protein and small molecule ligand.

Attention is a versatile mechanism and can also be applied to structural information

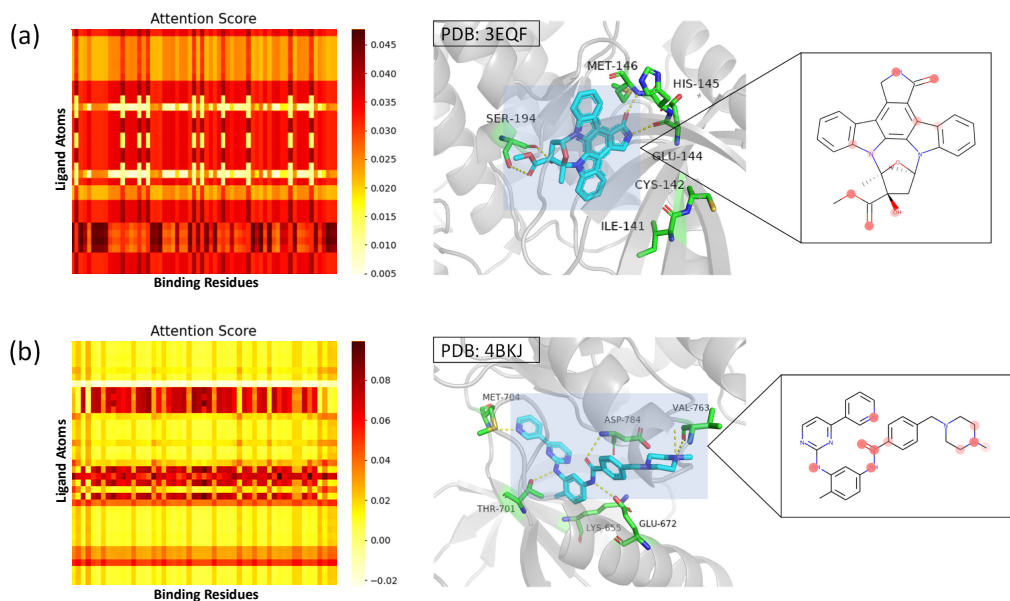


Figure 3: Sample Attention Weights for Relating Protein and Ligand. The heatmaps on the left help visualize the weighted importance of select protein residues and ligand atoms in a PLI. Structural views of the protein-ligand binding pocket are shown in the middle, with insets of the 2D ligand structures on the right. The colored residues and red color highlights indicate AAs in the protein binding pocket and ligand atoms with high attention scores. Adapted and modified from Figure 7 of Wu et al.¹¹⁰ used with permission under license CC BY 4.0.

such as the spatial coordinates of individual atoms or contact maps of protein-ligand complexes.^{111–113} The structural information of proteins and ligands can be well-represented by a graph with nodes representing AAs or atoms, and edges representing chemical bonds or amino acid contacts. Edges may also represent other predefined relationships or constraints between nodes. Integrating attention mechanisms into Graph Neural Networks (GNNs), a class of NNs specialized for processing graphs, has been increasingly-used for the study of PLIs.^{114–117} GNNs use "message-passing" whereby each node's embedding is updated iteratively based on information from connected nodes. Each connection can be assigned a weight that quantifies the likelihood of interdependence between connected nodes. For example, a cysteine residue may have a higher weight for a nearby cysteine than a nearby glycine due to the potential to form a disulfide bond between cysteines. GNNs are often augmented further, for example, by the addition of an attention mechanism to prioritize connected nodes during

message-passing. [111,112,115-118](#)

5.2.3. Transformers

While attention mechanisms have been quite beneficial for the predictive success of NLP methods, the "transformer" architecture pioneered in 2017 has also been instrumental in advancing these capabilities.⁹⁷ Transformers are a type of NN architecture that divides attention mechanisms into multiple parallel operations, each applying a different set of weights to the input data sequence. Several relationships between tokens are captured and processed simultaneously, dramatically improving the efficiency with which human text can be processed. The transformer architecture is the foundation of popular large language models such as ChatGPT²¹ and is a key component of DeepMind's AlphaFold system.^{15,119} Transformers have become widely used in bioinformatics, for DNA, RNA, and protein sequence analysis, as well as gene-based disease predictions and PLI predictions.¹²⁰

Transformers are designed to solve the problem of "sequence transduction" or the conversion of an input sequence of ordinal data into a predicted output sequence, such as a translated text or a vector representation.¹²¹ In NLP, this is called machine translation, whereby the input sequence, for example, could be a sentence in English and the output sequence is its French counterpart. Formally, the transformer architecture is an extension of the so-called "encoder-decoder" architecture, a state-of-the-art sequence-transduction method.⁹⁷ The premise of encoder-decoders is that sequentially ordered input data (e.g., English text, protein sequences, SMILES) can be "compressed" or encoded by a lower-dimensional fixed-length vector with minimal information loss. If the most "useful" or informative features of the data can be extracted and represented by this compressed or reduced vector representation, then the implicit rules/structures contained within the input data have been "encoded". Typically, in this reduced representation (called the "latent" space), inputs with similarly informative characteristics appear close to one another. These compressed vectors can subsequently be "decoded" or expanded to an output representation of choice

to complete the transduction task. These transduction tasks naturally align with the goal of text translation from one language to another. Importantly, transformers differ from traditional encoder-decoder models by incorporating the attention mechanism. Attention allows latent representations to vary in length, thus eliminating a fundamental constraint of encoder-decoder models: that every input sequence, regardless of length, be represented by a fixed-length vector in the latent space. Given their improved precision, transformers are widely used today (especially for long input sequences), especially given their inherent parallel architecture, which makes processing datasets with billions of items feasible. As compared to LSTMs, transformers are architecturally more complex and tend to achieve better performance.¹²²⁻¹²⁵ Even so, transformer performance is not always the best, particularly when dealing with small datasets on the order of thousands of items.¹²⁶⁻¹²⁸ In the biological domain, transformers have been applied to the prediction of protein-protein binding affinities,¹²⁹ post-translational modifications,¹³⁰ and quantum chemical properties of small molecules.¹³¹

Early applications of transformers for the study of PLIs involved simply retraining existing models designed for human language inputs;^{132,133} surprisingly, these transformers surpassed existing state-of-the-art models for predicting binding affinities.^{134,135} As new transformers were developed specifically to handle protein sequence data, predictive performance for PLIs improved.¹³⁶⁻¹³⁸ These developments included preemptively dividing the texts into subsequences to determine which amino acids contribute to binding and merging embeddings from different transformers to provide multiple representational perspectives. More recently, transformers have been adapted for use with other data types, such as protein structures and images, as well as for predicting PLI properties beyond binding affinity, e.g., binding poses.^{87,139}

5.3. Fusion of Protein-Ligand Representations: Concatenation or Cross-Attention

Once candidate interacting protein and ligand embeddings are extracted, they need to be fused for an interaction pattern to emerge. Method development for extracting embeddings from protein and ligand sequence data has been the primary focus to date such that approaches for fusion have been somewhat neglected, although they have garnered greater interest in recent years. A naive method for fusion is to simply end-to-end concatenate protein and ligand embedding vectors. A more refined approach could involve advanced data structures like graphs, whereby information such as coordinates of protein and ligand is used not only to build a graph representation but is also incorporated into an attention mechanism to account for local factors such as polarity or size.^{113,115,140} A mechanism of “cross-attention” could be incorporated into the fusion approach whereby the importance between the different *token representations* of the protein and ligand are directly calculated^{108,109,141} in an attempt to mirror the underlying interaction of a protein with a ligand.¹¹⁴ Cross-attention has been shown to be at least as competitive in predictive PLI tasks as other fusion methods,¹⁴⁰ and an improvement over the use of separate, independent attention mechanisms for both protein and ligand.¹⁴²

While fusion appears to be a natural and important component for NLP studies of PLIs, some models circumvent the idea of fusion altogether in lieu of considering only protein or ligand representations alone. For example, Wang et al.’s CSCConv2D algorithm only embeds ligand information.¹⁴³ An individual model is trained separately for each protein to predict that protein’s compatible ligands, resulting in the creation of hundreds of models. Although the study focuses on predicting PLIs, protein information is only incorporated indirectly by labeling ligands as either binding or non-binding to a given protein during its model’s training. Nonetheless, protein-only or ligand-only models are rare, with most contemporary NLP-PLI models considering both protein and ligand together through a fusion step.

Mixed-data approaches aimed at combining different data types for protein and/or ligand

(e.g., sequence + structure; sequence + image,¹³⁹ or both sequence and structure for protein + structure for ligand¹⁴¹) have further spurred studies into which input format is best for the protein or ligand. Mixed-data models may use a variety of architectures such as an LSTM or transformer for a protein sequence and a GNN for ligand structures.^{55,56} Combining multiple state-of-the-art embeddings for both sequence and structure has outperformed sequence-only baselines.⁵⁶ Despite the increased complexity involved in handling sequence and structural data simultaneously, mixed-data models are advantageous for both the ease-of-use of protein sequences and the completeness of ligand structural representations.

Although under-explored, combining multiple embeddings for each protein and ligand input in the fusion process may be beneficial. It has been suggested that different protein encoders for extracting features may gather different but relevant information to improve predictive outcomes.¹⁴⁴ In the algorithm, DeepPurpose, Huang et al. pursued a library approach that offered fifteen different protein and ligand embeddings (including transformer and RNN) to be combined and fed into a small NN to generate binary binding and/or continuous binding affinity predictions.¹⁴⁵ This menu-option approach enables users to compare feature extractors and find the best protein and ligand embeddings for their research. Another approach is to combine multiple embeddings through operations such as component-wise multiplication or component-wise difference, as each embedding could represent a different set of features.^{137,144} Shen et al.'s SVSBI algorithm¹³⁷ demonstrated how a higher-order embedding, by concatenating three different transformer embeddings, could outperform several state-of-the-art baselines in the prediction of binding affinity, including those based on individual transformers alone.

5.4 Prediction of Target Variables

Ultimately, specific research questions must motivate which relevant PLI target variables are to be predicted by the ML models constructed. These models often consist of one or more fully-connected layers with relatively few parameters than the NNs used for feature

extraction or fusion. The purpose of these layers is to utilize the latent protein and ligand features to predict an output target variable such as binding affinity or a binary indication of whether a pairing interacts. Thus, the fused protein and ligand embeddings are passed through these final layers to compute the prediction. Embeddings that effectively capture important underlying features can also be applied to predict other useful properties beyond binding affinity such as protein and ligand solubility.^{94,95}

Table 2: Sequence-Based PLI Prediction Models

Model Name	Extraction		Fusion	Prediction
	Protein Extractor	Ligand Extractor		
<i>LSTM</i>				
Affinity2Vec ¹⁰²	ProtVec	Seq2Seq	Heterogeneous Network	Gradient-Boosting Trees (R)
DeepLPI ¹⁰³	ResNet	ResNet	Concatenation with LSTM	FCN (C, R)
FusionDTA ¹⁰⁴	BiLSTM	BiLSTM	Concatenation with Linear Attention	FCN (R)
<i>Transformer</i>				
Shin et al. ¹³⁵	CNN	Transformer	Concatenation	FCN (R)
MolTrans ¹³⁶	Transformer	Transformer	Interaction Matrix ^a with CNN	FCN (C)
ELECTRA-DTA ¹³⁴	CNN with Squeeze-and-Excite Mechanism	CNN with Squeeze-and-Excite Mechanism	Concatenation	FCN (R)
MGPLI ¹³⁸	Transformer, CNN	Transformer, CNN	Concatenation	FCN (C)
SVSBI ¹³⁷	Transformer, LSTM, and AutoEncoder	Transformer, LSTM, and AutoEncoder	k-embedding fusion ^b	FCN, Gradient-Boosting Trees ^c (R)
<i>Non-Transformer Attention</i>				
DeepCDA ⁸⁶	CNN with LSTM	CNN with LSTM	Two-Sided Attention ^c	FCN (R)
HyperAttention-DTI ¹⁰⁹	CNN	CNN	Cross-Attention, Concatenation	FCN (C)
ICAN ¹⁰⁸	Various	Various	Cross-Attention, Concatenation	1D CNN (C)
<i>Other NLP Methods</i>				
GANsDTA ¹⁴⁶	GAN Discriminator	GAN Discriminator	Concatenation	1D CNN (R)
Multi-PLI ¹⁴⁷	CNN	CNN	Concatenation	FCN (C, R)
ChemBoost ⁸⁹	Various	SMILESVec	Concatenation	Gradient-Boosting Trees (R)

Note: A model’s task of Classification (C) and/or Regression (R) is denoted beside the "Prediction" column entries in parenthesis. Definitions for specific terms may be found in the Glossary (Table 6).

Terms Defined by the Cited Authors: ^a**Interaction Matrix:** Output from dot product operations to measure interactions between protein sub-sequence and ligand sub-structure pairs. ^b**k-embedding fusion:** The use of machine learning to find an optimal combination of lower-order embeddings via different integrating operations. ^c**Two-sided Attention:** Attention mechanism that computes scores using the products of both pairs of protein/ligand fragments and protein/ligand feature vectors.

5.5. Evaluation

Evaluation is typically performed by comparing statistical metrics between models on the same test datasets. Evaluation metrics vary by task: classification predictions can be as-

Table 3: Structure-Based PLI Prediction Models

Model Name	Extraction		Fusion	Prediction
	Protein Extractor	Ligand Extractor		
<i>Transformer</i>				
UniMol ⁸⁷	Transformer-Based Encoder	Transformer-Based Encoder	Concatenation	Transformer-Based Decoder (R)
<i>Other Attention</i>				
Lim et al. ¹¹⁸	GNN	GNN	Attention	FCN (C)
Jiang et al. ¹¹¹	GCN	GCN	Concatenation	FCN (R)
GEFA ¹¹²	GCN	GCN	Concatenation	FCN (R)
Knutson et al. ¹¹⁴	GAT	GAT	Concatenation	FCN (C, R)
AttentionSite-DTI ¹¹⁷	GCN with Attention	GCN with Attention	Concatenation, Self-Attention	FCN (C, R)
HAC-Net ¹¹⁵	GCN with Attention Aggregation	GCN with Attention	Combined Graph Representation	FCN (R)
BindingSite-AugmentedDTI ¹¹⁶	GCN with Attention	GCN with Attention	Concatenation, Self-Attention	Various (R)
PBCNet ¹¹³	GCN	Message-Passing NN	Attention	FCN (R)

Note: A model’s task of Classification (C) and/or Regression (R) is denoted beside the "Prediction" column entries in parenthesis. Definitions for specific terms may be found in the Glossary (Table 6).

sessed via precision, recall, and F1 score metrics whereas regression predictions are often evaluated relative to the ground-truth test data via concordance index and mean square error metrics.^{155,156} Pre-made datasets such as PDBBind⁵⁹ frequently come with both training and test datasets to enable fair comparisons with other established models. Models aiming to be generalizable across several types of PLIs should ideally be evaluated on several different sets of proteins and ligands.

While ML models can be assessed through the aforementioned statistical metrics, the practical utility of PLI predictive models and their predictive accuracy in real-world cases is best validated by PLI domain experts in the field.¹⁵⁷ For example, if a model is designed to predict binding affinities, a set of predictions generated *in silico* would be best confirmed through *in vitro* experimentation. This can serve two purposes: justifying a model’s use where it can be most effective and creating an opportunity for future interdisciplinary collaboration between ML practitioners and PLI domain experts in computational and experimental biology.

Table 4: Mixed Representation PLI Prediction Models

Model Name	Input Type	Extraction		Fusion	Prediction
		Protein	Ligand		
<i>LSTM</i>					
Zheng et al. ¹⁴⁸	P: Struct. L: Seq.	Dynamic CNN ^a with Attention	BiLSTM with At- tention	Concatenation	FCN (C)
DeepGLSTM ⁵⁵	P: Seq. L: Struct.	BiLSTM with FCN	GCN	Concatenation	FCN (R)
<i>Transformer</i>					
Transformer-CPI ⁵⁶	P: Seq. L: Struct.	Transformer En- coder	GCN	Transformer Decoder	FCN (C)
DeepPurpose ¹⁴⁵	P: Seq. L: Either	4 Various Encoders	5 Various Encoders	Concatenation	FCN (C, R)
CAT-CPI ¹³⁹	P: Seq. L: Image	Transformer En- coder	Transformer En- coder	Concatenation	CNN and FCN (C)
<i>Non-Transformer Attention</i>					
Tsubaki et al. ¹⁴⁹	P: Seq. L: Struct.	CNN	GNN	Attention and Concatenation	FCN (C)
DeepAffinity ¹⁵⁰	P: Seq. L: Struct.	RNN-CNN with Attention	RNN-CNN with Attention	Concatenation	FCN (R)
MONN ¹⁵¹	P: Seq. L: Struct.	CNN	GCN	Pairwise Inter- action Matrix ^b , Attention	Linear Regression (C, R)
GraphDTA ¹⁴⁰	P: Seq. L: Struct.	CNN	4 GNN Variants	Concatenation	FCN (R)
CPGL ¹⁵²	P: Seq. L: Struct.	LSTM	GAT with Attention	Two-Sided Attention ^c , Concatenation	Logistic Regres- sion (C)
CAPLA ¹⁴¹	P: Both L: Struct.	Dilated Convolu- tional Block	Dilated Convolu- tional Block with Cross-Attention to Binding Pocket	Cross-Attention, Concatenation	FCN (R)

Note: A model’s task of Classification (C) and/or Regression (R) is denoted beside the "Prediction" column entries in parenthesis. Definitions for specific terms may be found in the Glossary (Table 6). The input representations for sequence and structure are abbreviated for brevity.

Terms Defined by the Cited Authors: ^a**Dynamic CNN:** ResNet-based CNN modified to handle inputs of variable lengths by padding the sides of the input with zeroes. ^b**Pairwise Interaction Matrix:** A [number of atoms]-by-[number of residues] matrix in which each element is a binary value indicating if the corresponding atom-residue pair has an interaction. ¹⁵¹ ^c**Two-sided Attention:** Attention mechanism that uses dot product operations between protein AA and ligand atom pairs, while taking matrices of learned weights into account.

6. Challenges and Future Directions

Generative AI and NLP techniques have revolutionized how we tackle tasks related to human language. Early successes of NLP methods in discerning the "rules" of protein structure (as exemplified by AlphaFold¹⁵) suggest significant potential for NLP to transform our approach to studying PLIs. While many innovations in the NLP computational toolkit for PLIs have emerged in recent years, several practical hurdles remain, limiting the impact and potential insights derivable from the ML approaches. This section presents an overview of the many challenges confronting the PLI field and suggests various avenues to address them.

Table 5: Sequence-Structure-Plus PLI Prediction Models

Model Name	Extraction			Fusion	Prediction
	Protein Extractor	Ligand Extractor	Additional Features Used		
<i>LSTM</i>					
HGDTI ¹⁵³	BiLSTM	BiLSTM	Disease and Side Effect Information	Concatenation	FCN (C)
ResBiGAAT ⁵⁷	Bidirectional GRU with Attention	Bidirectional GRU with Attention	Global Protein Features	Concatenation	FCN (R)
<i>Transformer</i>					
Gaspar et al. ⁹⁰	Transformer or LSTM	ECFC4 Fingerprints	Multiple Sequence Alignment Information	Concatenation	Random Forest (C)
HoTS ¹⁵⁴	CNN	FCN	Binding Region	Transformer Block	FCN (C, R)
PLA-MoRe ⁵⁸	Transformer	GIN and AutoEncoder	Bioactive Properties	Concatenation	FCN (R)
AlphaFold 3 ¹⁶	Attention-Based Encoder ^a	Attention-Based Encoder ^a	Post-Translational Modifications, Multiple Sequence Alignment Information	Attention	Diffusion Transformer ^b
<i>Other NLP Methods</i>					
MultiDTI ⁸⁸	CNN with FCN	CNN with FCN	Disease and Side Effect Information	Heterogeneous Network	FCN (C)

Note: A model’s task of Classification (C) and/or Regression (R) is denoted beside the "Prediction" column entries in parenthesis. Definitions for specific terms may be found in the Glossary (Table 6).

Terms Defined by the Cited Authors: ^a***Atom Attention Encoder:** An attention-based encoder that uses cross-attention to capture local atom features. ^b***Diffusion Transformer:** A transformer-based model that aims to remove noise from predicted atomic coordinates until a suitable final structure is output.

6.1. Lack of "True Negatives"

A common challenge in today’s data-driven ML paradigm is the limited availability of abundant, high-quality, and labeled data. In PLI studies, there is a particular lack of bona fide "negative examples", i.e., data for ligand-like molecules that do not bind a protein of interest that are critical for model training. In "supervised" ML,¹⁵⁸ models are trained on data with labels of whether a protein-ligand pair is binding or non-binding, and protein-ligand data spanning the full spectrum of interaction/no-interaction are necessary for models to 'learn'. When a similar situation is encountered in other ML problems, a common approach is to select random data points not explicitly labeled as "positive" and to declare them as "negative". This would be equivalent to assigning random ligands to each protein and treating them as negative PLI examples. Unfortunately, given the complexity and specificity of PLIs, these are often *trivial* negative examples since molecules that do not interact with a

protein of interest *and* are dissimilar to the "true" ligands embody little information from which ML models can learn. Manually curating protein-ligand pairs that display weak interactions or lower binding affinities is an option for addressing this problem, although this is time-consuming and labor-intensive.

Acquiring the requisite negative data for classification studies is tied to experimental studies that conclusively determine whether pairings bind. The availability of informative negative PLI data requires deliberate efforts of domain experts who recognize the importance of generating, curating, and reporting such data, which are rarely publicized or emphasized in the literature regardless of data type.¹⁵⁹⁻¹⁶¹ Learning from positive data only or with unlabeled data is, therefore, an active field of study, with many attempts applying "unsupervised" or "semi-supervised" methods¹⁶² (see^{146,163} for examples related to PLI prediction). Compared with supervised models, un-/semi-supervised models typically require larger datasets of tens to hundreds of thousands of PLIs. Furthermore, the associated network architectures may be more computationally intensive.¹⁴⁶ In cases where negative data does exist but at a significantly reduced quantity, additional remedies may be attempted. For example, classification studies of PLIs can adjust the distribution of ligands to ensure *equal proportions* in the positive and negative examples represented; this has been shown to mitigate the issue of an overabundance of positive data.¹⁶⁴ Future studies should resolve the lack of readily available non-interacting protein-ligand pairs, perhaps through mining the scientific literature for meaningful non-binding pairs.

6.2. Diversity Bias in PLI Datasets

Many PLI datasets possess underlying bias with respect to either the diversity or types of proteins and ligands represented, which hinders the effectiveness of ML algorithms. Training with *insufficiently different* data points can lead to poor generalizability and predictive performance when extended to real-world examples not represented in the training dataset. For example, binding affinity predictors trained on the popular PDBBind dataset⁵⁹ with

both protein and ligand information represented performed no better than those trained on protein- or ligand-only information.¹⁶⁵ This implies that the PDDBind dataset is biased and that the protein-ligand trained model failed to discern the mechanics of binding and rather "memorized" the most popular representatives or non-informative patterns within the dataset. The commonly used DUD-E⁶² dataset of bioactive compounds and respective protein targets was found to suffer from a similar problem: PLI binding classification models that differentiated binders/non-binders to a high degree of accuracy resulted only because the binders and non-binders were of different shape classes and not because they embedded any relevant information about the protein-ligand interface.^{165,166} The existing literature suggests that this is a problem of quality over quantity, as memorization-related biases in PLI models are *not* alleviated by merely increasing the dataset size or removing over-represented items.¹⁶⁷ The presence of bias is understandable, given how idiosyncratic research interests in biological or pharmaceutical fields shape the particular proteins and subsets of ligands that are studied and the type of PLI data generated and made available. Given that models trained on biased data often fail in practical, real-world prediction tasks, the creation of high-quality, well-balanced, and unbiased PLI datasets is essential to the future of ML-based PLI studies and should be made a priority.

To train more generalizable models, systematic datasets with proteins and ligands beyond those of biological and pharmaceutical interest need to be evaluated. One way around the experimental challenges of generating sufficient protein-ligand data may be through high-throughput molecular dynamics simulations and/or docking studies of PLIs using AlphaFold-predicted¹⁵ protein structures. Although current simulation methods are time intensive, the availability of powerful computing clusters and trends towards increasingly powerful GPU hardware may make this approach feasible in the not-too-distant future, and the benefits may be worth investing in this pursuit. This approach could be automated, requiring far less human intervention than laboratory experiments, and can yield valuable binding pocket information for better structure-based ML predictions.

6.3. The Limitations of "Language-ness" in Protein and Ligand Text Representations

As compared to human languages, both proteins and ligands have significant structural and ontological differences that have to be accounted for when designing a model. The following nuances have driven investigation into modifying existing NLP architectures to accommodate for protein and ligand representations.

In linguistics, a "word" is a complete unit of meaning that a reader can recognize. However, for protein sequences, such corresponding units are not easily demarcated. It would be dubious to assume AAs are equivalent to "words" because the roles of individual AAs are highly dependent on their context and environment. The meaning of a word in a human language may be independent of its surroundings if the word has only one definition; however, an amino acid has "meaning" across several levels/dimensions, influencing secondary structure, tertiary structure, motif function, and/or binding interactions. Conversely, it is also difficult to view motifs or domains as "words" since not all regions of a protein are independent of one another.¹⁶⁸ This lack of word-equivalence is one motivation behind "sub-word" tokenization methods that attempt to discern a hierarchy of word equivalents in protein sequences.⁸⁵ Protein sequences are also different from human languages in the length scale of interactions and the number of long-distance interactions that contribute to a 3D structure. While human language texts have distant relationships, such as between the subject and pronoun of a sentence or a passage that foreshadows another in an essay, these relations can be deduced by a reader and remain relatively sparse on a per-sentence basis. In contrast, AAs have numerous distant relationships and cannot be thoroughly predicted by even an expert in protein biochemistry without the assistance of computational or experimental tools, thus adding a layer of complexity to the analysis of AA sequences that can be well compensated for by ML approaches.

In the chemical space, the SMILES format is dissimilar from human languages in the large variation in text length and in the difficulty of finding an ideal tokenization scheme. First,

the lengths of SMILES strings could vary even more than those of protein sequences, ranging from listing each atom of small molecules to those constituting entire proteins, although protein-protein interactions are generally considered a separate problem. The SMILES format is less practical to use for larger molecules, since structural graphs can provide a more compact and accurate representation of atoms in a large three-dimensional structure. A further disadvantage of using SMILES is that it is difficult to intuitively discern “word” equivalents within the string. Individual branches separated by parenthesis could be viewed as words,¹⁶⁹ but this is only practical for small branching groups. Moreover, the handling of nesting parentheses in SMILES molecules can be problematic and has become a major limiting factor in ML models designed to generate novel molecules.¹⁷⁰ The sum of these SMILES shortcomings has led to the development of alternative chemical representations for computational studies such as DeepSMILES and SELFIES.^{171,172} Although promising, these alternate forms have rarely been used in PLI studies to date. The question remains whether a three-dimensional molecule can be truly mapped to a text representation in a way that preserves all relevant structural information for PLIs.

6.4. Interpretable Design in PLI Predictions

Catalyzed by the open-data movement and widely accessible machine-learning tools, hidden or explicit patterns are discovered from datasets through weighted mathematical operations that are difficult to interpret and yet many effective predictive models have been developed. A majority of ML studies fail to consider designing human-friendly interpretations of *how* their models’ predictions are calculated. This is a significant contemporary challenge that has prompted the growth of explainable AI (XAI) as an active research field.

To build trust among biologists and broaden scientific acceptance, future ML models must be more understandable to end-user biologists than provided by common “black-box” designs.³⁵ One potential approach for bridging the “explainability” gap in PLI studies is the use of attention weights to corroborate existing protein-ligand contacts (cf. Fig. 3).^{56,86,104,108,109}

Attention weights provide a degree of interpretability by highlighting binding regions in PLI models that converge with higher weight values. Given the reality of "false positives" whereby higher binding weights are inadvertently assigned to non-binding regions, attention weights alone may not constitute a satisfactory basis for explaining or inferring what regions govern binding interactions. Unfortunately, a systematic assessment of 'false positives' in attention weights has yet to be performed, leaving it unclear whether they are a reliable metric.¹⁵⁶ Such potential inaccuracies are but one facet of a larger debate on whether attention weights provide sufficient explanatory power for PLI models.¹⁷³

While NLP presents attention mechanisms as one possible avenue, other methods of explainability have also started to be applied to the study of PLIs. One example is the use of a game-theoretical approach to compute "Shapley values", which quantify the importance of individual features by evaluating each feature's contribution to the final prediction across all possible combinations of those features.^{174,175} Visualization may also be a great tool for identifying possible binding interactions. For example, graph visualization can help depict the bonds between an interacting protein and ligand, and "saliency maps" can highlight specific subregions of protein and ligand that are the most influential in the model's prediction by determining how changes in individual input features affect the output.¹⁷⁶ Several underutilized avenues for establishing interpretability remain,¹⁷⁷ but none have been established as a state-of-the-art; determining a standard method of interpretability for PLI prediction models will be critical for the field.

6.5. The Insufficiency of an NLP-only Approach for PLI Studies?

While NLP offers beneficial strategies for the study of PLIs, it is far from a panacea and is often complemented by insights from other disciplines. Existing attention-based, state-of-the-art NLP models are limited by the need for substantial amounts of training data for the best results.¹⁷⁸ There may be opportunities for other disciplines of computer science to contribute positively to the PLI field. For example, it may be more fruitful to incorpo-

rate computer vision techniques that are better at handling structural information over NLP techniques that are designed for handling text.¹⁷⁹

Several studies have combined NLP with more unusual architectures or complementary approaches. For example, Zhao et al. created an algorithm that uses so-called generative adversarial networks (GANs) as a means of embedding protein sequences and SMILES independently.¹⁴⁶ GANs feature a dual NN architecture: a generator that creates artificial data points and a discriminator that is trained to distinguish between real and artificial data. Both components were trained together in a process akin to an evolutionary arms-race, as the discriminator repeatedly learns to identify key features of the input that help distinguish real data from increasingly realistic artificial data. Zhao et al. demonstrated competitive results relative to selected benchmarks even though the efficacy of their GAN was stated to be limited by the small dataset used but would likely perform better if trained on at least thousands of diverse proteins and ligands.¹⁴⁶ In other studies, computer vision methods have been used to combine images of proteins or ligands as inputs alongside their text representations; features across the two modalities enable attention mechanisms to capture cross-feature correlations across data types.^{139,148} While the success of the aforementioned hybrid strategies did not exceed the performance of other neural networks in PLI predictive tasks,¹⁸⁰ they demonstrate the potential for innovation using advances from other sub-domains of ML and computer science beyond NLP.

Biological domain knowledge is crucial both for framing the computational challenges related to ML and for identifying best practices for handling protein and ligand data. Approaches that are grounded in a deep understanding of the underlying domain-specific science have proven to be forerunners in the practical success of ML methods, as demonstrated by the AlphaFold initiative’s sophisticated use of sequence evolutionary information.^{15,119} The study of PLIs may eventually outgrow NLP methods, but for the foreseeable future, advances in NLP will continue to have a significant impact. Collaborations between experts in biological and computational domains will be critical for catalyzing further innovations in

what is an interdisciplinary goal.

7. Conclusion

Natural language processing (NLP), a sub-discipline of machine learning (ML), offers new tools for both experimental and computational researchers to accelerate exploratory studies in structural biology. The prediction of protein-ligand interactions (PLIs) can be re-imagined through NLP by treating protein and ligand representations like text. Protein sequences resemble readable text with inherent meaning to be inferred, while the SMILES format for chemical compounds allows limited NLP application to small molecules. Current efforts seek to leverage multiple or augmented SMILES representations to address these limitations.

Approaches to tackling PLI prediction tasks using sequence-only data, structural data, or a combination of both, have all yielded successful predictions, although the advantage of one input data type over others remains unclear. Sequence-only data approaches are simple and amenable to NLP but requires a significant abstraction of chemical information; structural data is informationally rich but computationally expensive to handle, while combining both sequence and structural data types offers balance at the expense of complexity.

The transformer architecture, in general, and attention mechanisms, in particular, have yielded the most promising NLP-based PLI prediction results to date. Incorporating complementary data (e.g., multiple sequence alignments, ligand polarities, etc.) can improve predictive success but at a significant increase in computational cost. After data selection and preparation, all methods have followed a general ML Extract-Fuse-Predict model creation framework of: (i) extracting feature embeddings for protein and ligand, (ii) fusing protein and ligand embeddings, and (iii) making predictions based on the created ML model.

The first step of dataset selection is crucial for any ML-based study of PLIs, and no single dataset can satisfy all needs, with many suffering from missing data or lack of negative data. Datasets must align with specific research goals, requiring thoughtful consideration as

to what inputs, formats, and target variable(s) are selected for the ML model. Appropriate tokenization and embedding methods, which convert proteins and ligands into numerical representations, are vital for a successful model. Atoms or amino acids typically serve as tokens, and neural networks (NNs) have helped to identify hidden patterns more quickly. NLP-inspired NNs, such as Long Short-Term Memory NNs, along with attention mechanisms and transformer architectures, have shown particular promise for understanding PLIs. A modular approach combining multiple embeddings can capture diverse perspectives, improving prediction accuracy, especially for the prediction of binding affinities. After appropriate embeddings are obtained, graph-based methods and cross-attention mechanisms have been shown to be effective in combining data from diverse sources.

NLP has been central to ML studies of PLIs and has yielded promising results, although many challenges remain. Explaining ML model predictions is essential for their trustworthiness and acceptance. Current explanatory metrics, such as attention weights and Shapley values, offer some degree of interpretability but remain to be fully validated. A major challenge is the lack of well-annotated non-binding protein-ligand pairs, or "negative data". Unsupervised methods or manually curated selections of non-binding pairs are potential solutions. Popular PLI datasets may contain biases that cause models to "memorize" idiosyncratic patterns rather than "learn" the true mechanics of PLIs. Ensuring balanced training datasets (positive vs. negative data, number of proteins vs. ligands, etc.) would be essential to avoid such bias.

As protein and ligand sequence representations differ from human language, it may be difficult to capture their complexity with NLP methods alone, especially as much of the variation in protein function can often be explained by simple amino acid interactions rather than complex higher-order interactions.¹⁸¹ While NLP has contributed significantly to advancing our study of PLIs, future improvements may come from both modifying machine learning architectures and incorporating nuanced biological domain knowledge. For instance, the researchers behind AlphaFold-Multimer's protein-protein interaction prediction algorithm¹⁸²

created an interface-aware protocol that crops protein structures to reduce computational burden and the representation of non-interfacial amino acids while maintaining an important balance of interacting and non-interacting regions. Some researchers have also integrated mass spectrometry data to improve model predictions of protein complexes.¹⁸³ More recently in AlphaFold3,¹⁶ a diffusion layer has been added to Alphafold’s previous workflow to enable the study of PLIs. Time will tell to what degree AlphaFold3 will advance predictions of PLIs but progress in PLI research will undoubtedly require interdisciplinary collaborations between computer scientists, chemists, and biologists.

Although it is best practice to evaluate model performance against ground-truth experimental results or results from physics-based computer simulations, few studies to date have benchmarked their model predictions in this way. Formal competition may prove to be a promising avenue for future advances in PLI prediction. Other grand challenges, such as protein folding and protein assembly, have had significant progress facilitated through competitions like Critical Assessment of Structural Prediction (CASP)¹⁸⁴ and Critical Assessment of Prediction of Interactions (CAPRI).^{185,186} These well-adjudicated competitions use unpublished test sets for objective model comparisons. Milestone algorithms like AlphaFold¹⁸⁷ and RosettaFold¹⁸⁸ were formed, improved, and refined through the crucible of such contests. Creating a dedicated competition devoted to protein-ligand interactions could similarly inspire innovation and catalyze seminal algorithmic advances for PLI prediction.

Acknowledgement

This work was supported in part by NIGMS/NIH Institutional Development Award (IDeA) #P20GM130460 to J.L, NSF award #1846376 to E.F.Y.H, and University of Mississippi Data Science/AI Research Seed Grant award #SB3002 IDS RSG-03 to J.M., J.L., T.L, and E.F.Y.H.

Table 6: Glossary of Terms That Appear in the Tables

Term	Definition
AutoEncoder	A neural network tasked with compressing and reconstructing input data, often used for feature learning. ¹⁸⁹
Dilated Convolutional Block	Convolutional Neural Network operations with defined gaps between kernels, which can capture larger receptive fields with fewer parameters.
ECFC4 Fingerprint	A molecular fingerprint that encodes information about the presence of specific substructures within a diameter of 4 bonds from each atom. ¹⁹⁰
FCN	Fully-Connected Network, a feedforward Neural Network where each neuron in one layer connects to every layer in the next. FCNs can also be referred to as Multi-Layer Perceptrons.
GAN Discriminator	An NN part of Generative Adversarial Networks (GAN) that learns important features to distinguish between real and artificial data.
GAT	Graph Attention Network, a type of Graph Neural Network that uses attention mechanisms to deciding the value of neighboring nodes to a given node when updating a node’s information. ¹⁹¹
GCN	Graph Convolutional Network, a type of Graph Neural Network that aggregates neighboring node features through a first-order approximation on a local filter of the graph. ¹⁹²
GIN	Graph Isomorphism Network, a type of Graph Neural Network that uses a series of functions to ensure embeddings are the same no matter what order nodes are presented in. ¹⁹³
Gradient-Boosting Trees	A machine learning technique where many decision trees are trained in order, such that the next tree learns from the misclassified samples of the previous tree. All trees are then used to "vote" on results of each input.
GRU	Gated Recurrent Unit, a simplified version of Long Short-Term Memory that similarly uses a gating mechanism to retain and forget information, but is less complex than Long Short-Term Memory. ¹⁹⁴
Heterogeneous Network	A graph where nodes and edges represent different types of information, often used to convey complex relationships in biological systems (e.g. drug, target, side-effect, etc.).
Message-Passing NN	Type of Graph Neural Network that computes individual messages to be passed between nodes so that representations for each node contain information from its neighbors. ¹⁹⁵
ProtVec	A method for representing protein sequences as dense vectors using skip-gram neural networks. ¹⁹⁶
Random Forest	A machine learning method where many decision trees are constructed, and the result of the ensemble is the mode of the individual tree predictions.
ResNet	Short for Residual Network. A neural network architecture that speeds up training by learning functions to substitute for layer operations, allowing for the "skipping" of layers and faster training. ¹⁹⁷
Seq2Seq	A machine learning method used for language translation in NLP, featuring an encoder-decoder structure. ¹⁹³
SMILESVec	Previous work from authors. 8-character ligand SMILES fragments are assigned a vector through a single-layer neural network, and an input SMILES string’s vector is equal to the mean of fragment vectors present in that input SMILES. ¹⁹⁸
Squeeze-And-Excite Mechanism	Mechanism for Convolutional Neural Networks that uses global information to adapt the model to emphasize more important features. ¹⁹⁹

References

- (1) Sarkar, D.; Saha, S. Machine-learning techniques for the prediction of protein-protein interactions. *J. Biosci.* **2019**, *44*.
- (2) Huang, S.-Y.; Zou, X. Advances and challenges in protein-ligand docking. *Int. J. Mol. Sci.* **2010**, *11*, 3016–3034.
- (3) Chaires, J. B. Calorimetry and thermodynamics in drug design. *Annu. Rev. Biophys.* **2008**, *37*, 135–151.
- (4) Woodley, J. M. Protein engineering of enzymes for process applications. *Curr. Opin. Chem. Biol.* **2013**, *17*, 310–316.
- (5) Barbosa, A. J. M.; Oliveira, A. R.; Roque, A. C. A. Protein- and Peptide-Based Biosensors in Artificial Olfaction. *Trends Biotechnol.* **2018**, *36*, 1244–1258.
- (6) Vajda, S.; Guarnieri, F. Characterization of protein-ligand interaction sites using experimental and computational methods. *Curr. Opin. Drug Discov. Devel.* **2006**, *9*, 354–362.
- (7) Sousa, S. F.; Ribeiro, A. J. M.; Coimbra, J. T. S.; Neves, R. P. P.; Martins, S. A.; Moorthy, N. S. H. N.; Fernandes, P. A.; Ramos, M. J. Protein-Ligand Docking in the New Millennium A Retrospective of 10 Years in the Field. *Curr. Med. Chem.* **2013**, *20*, 2296–2314.
- (8) Morris, C. J.; Corte, D. D. Using molecular docking and molecular dynamics to investigate protein-ligand interactions. *Mod. Phys. Lett. B* **2021**, *35*, 2130002.
- (9) Lecina, D.; Gilabert, J. F.; Guallar, V. Adaptive simulations, towards interactive protein-ligand modeling. *Sci. Rep.* **2017**, *7*, 8466.

- (10) Ikebata, H.; Hongo, K.; Isomura, T.; Maezono, R.; Yoshida, R. Bayesian molecular design with a chemical language model. *J. Comput. Aided Mol. Des.* **2017**, *31*, 379–391.
- (11) Rives, A.; Meier, J.; Sercu, T.; Goyal, S.; Lin, Z.; Liu, J.; Guo, D.; Ott, M.; Zitnick, C. L.; Ma, J.; Fergus, R. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proc. Natl. Acad. Sci. U. S. A.* **2021**, *118*.
- (12) Cao, Y.; Shen, Y. TALE: Transformer-based protein function Annotation with joint sequence–Label Embedding. *Bioinformatics* **2021**, *37*, 2825–2833.
- (13) Wang, S.; Guo, Y.; Wang, Y.; Sun, H.; Huang, J. SMILES-BERT: Large Scale Unsupervised Pre-Training for Molecular Property Prediction. Proceedings of the 10th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics. New York, NY, USA, 2019; pp 429–436.
- (14) Chithrananda, S.; Grand, G.; Ramsundar, B. ChemBERTa: Large-Scale Self-Supervised Pretraining for Molecular Property Prediction. *arXiv [cs.LG]* **2020**,
- (15) Jumper, J. et al. Highly accurate protein structure prediction with AlphaFold. *Nature* **2021**, *596*, 583–589.
- (16) Abramson, J. et al. Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Nature* **2024**,
- (17) Zhou, M.; Duan, N.; Liu, S.; Shum, H.-Y. Progress in neural NLP: Modeling, learning, and reasoning. *Engineering (Beijing)* **2020**, *6*, 275–290.
- (18) Patwardhan, N.; Marrone, S.; Sansone, C. Transformers in the real world: A survey on NLP applications. *Inf.* **2023**, *14*, 242.

- (19) Bijral, R. K.; Singh, I.; Manhas, J.; Sharma, V. Exploring Artificial Intelligence in Drug Discovery: A Comprehensive Review. *Arch. Comput. Methods Eng.* **2022**, *29*, 2513–2529.
- (20) Ray, P. P. ChatGPT: A comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope. *Internet of Things and Cyber-Physical Systems* **2023**, *3*, 121–154.
- (21) Radford, A.; Narasimhan, K.; Salimans, T.; Sutskever, I. Improving language understanding by generative pre-training. <https://www.mikecaptain.com/resources/pdf/GPT-1.pdf>, Accessed: 2023-10-27.
- (22) Goodside, R., Papay, Meet Claude: Anthropic’s Rival to ChatGPT. <https://scale.com/blog/chatgpt-vs-claude>, 2023; Accessed: –.
- (23) Bing Copilot. [BingCopilot;https://copilot.microsoft.com/](https://copilot.microsoft.com/).
- (24) Rahul; Adhikari, S.; Monika NLP based Machine Learning Approaches for Text Summarization. 2020 Fourth International Conference on Computing Methodologies and Communication (ICCMC). 2020; pp 535–538.
- (25) Nasukawa, T.; Yi, J. Sentiment analysis: capturing favorability using natural language processing. Proceedings of the 2nd international conference on Knowledge capture. New York, NY, USA, 2003; pp 70–77.
- (26) Lample, G.; Charton, F. Deep Learning for Symbolic Mathematics. *arXiv [cs.SC]* **2019**,
- (27) Feng, Z.; Guo, D.; Tang, D.; Duan, N.; Feng, X.; Gong, M.; Shou, L.; Qin, B.; Liu, T.; Jiang, D.; Zhou, M. CodeBERT: A Pre-Trained Model for Programming and Natural Languages. *arXiv [cs.CL]* **2020**,

- (28) Ashok, V. G.; Feng, S.; Choi, Y. Success with style: Using writing style to predict the success of novelsd.
- (29) Barberá, P.; Boydston, A. E.; Linn, S.; McMahon, R.; Nagler, J. Automated text classification of news articles: A practical guide. *Polit. Anal.* **2021**, *29*, 19–42.
- (30) Wang, H.; Wu, H.; He, Z.; Huang, L.; Church, K. W. Progress in machine translation. *Engineering (Beijing)* **2022**, *18*, 143–153.
- (31) Bhasuran, B.; Natarajan, J. Automatic extraction of gene-disease associations from literature using joint ensemble learning. *PLoS One* **2018**, *13*, e0200699.
- (32) Pang, M.; Su, K.; Li, M. Leveraging information in spatial transcriptomics to predict super-resolution gene expression from histology images in tumors. *bioRxiv* **2021**, 2021.11.28.470212.
- (33) Bouatta, N.; Sorger, P.; AlQuraishi, M. Protein structure prediction by AlphaFold2: are attention and symmetries all you need? *Acta Crystallogr D Struct Biol* **2021**, *77*, 982–991.
- (34) Skolnick, J.; Gao, M.; Zhou, H.; Singh, S. AlphaFold 2: Why It Works and Its Implications for Understanding the Relationships of Protein Sequence, Structure, and Function. *J. Chem. Inf. Model.* **2021**, *61*, 4827–4831.
- (35) Adadi, A.; Berrada, M. Peeking inside the black-box: a survey on explainable artificial intelligence (XAI). *IEEE access* **2018**,
- (36) Box, G. E. P. Science and Statistics. *J. Am. Stat. Assoc.* **1976**, *71*, 791–799.
- (37) Geirhos, R.; Jacobsen, J.-H.; Michaelis, C.; Zemel, R.; Brendel, W.; Bethge, M.; Wichmann, F. A. Shortcut learning in deep neural networks. *Nature Machine Intelligence* **2020**, *2*, 665–673.

- (38) Outeiral, C.; Nissley, D. A.; Deane, C. M. Current structure predictors are not learning the physics of protein folding. *Bioinformatics* **2022**, *38*, 1881–1887.
- (39) Ferruz, N.; Höcker, B. Controllable protein design with language models. *Nature Machine Intelligence* **2022**, *4*, 521–532.
- (40) Ofer, D.; Brandes, N.; Linial, M. The language of proteins: NLP, machine learning & protein sequences. *Comput. Struct. Biotechnol. J.* **2021**, *19*, 1750–1758.
- (41) Ptitsyn, O. B. How does protein synthesis give rise to the 3D-structure? *FEBS Lett.* **1991**, *285*, 176–181.
- (42) Yu, L.; Tanwar, D. K.; Penha, E. D. S.; Wolf, Y. I.; Koonin, E. V.; Basu, M. K. Grammar of protein domain architectures. *Proc. Natl. Acad. Sci. U. S. A.* **2019**, *116*, 3636–3645.
- (43) Petsko, G. A.; Ringe, D. *Protein Structure and Function*; Primers in biology; Blackwell Publishing: London, England, 2003.
- (44) Shenoy, S. R.; Jayaram, B. Proteins: sequence to structure and function—current status. *Curr. Protein Pept. Sci.* **2010**, *11*, 498–514.
- (45) Garfield, E. Chemico-linguistics: computer translation of chemical nomenclature. *Nature* **1961**, *192*, 192.
- (46) Wigh, D. S.; Goodman, J. M.; Lapkin, A. A. A review of molecular representation in the age of machine learning. *Wiley Interdiscip. Rev. Comput. Mol. Sci.* **2022**, *12*.
- (47) Weininger, D. SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *J. Chem. Inf. Comput. Sci.* **1988**, *28*, 31–36.
- (48) Wang, Y.; Xiao, J.; Suzek, T. O.; Zhang, J.; Wang, J.; Bryant, S. H. PubChem: a public information system for analyzing bioactivities of small molecules. *Nucleic Acids Res.* **2009**, *37*, W623–33.

- (49) Degtyarenko, K.; de Matos, P.; Ennis, M.; Hastings, J.; Zbinden, M.; McNaught, A.; Alcántara, R.; Darsow, M.; Guedj, M.; Ashburner, M. ChEBI: a database and ontology for chemical entities of biological interest. *Nucleic Acids Res.* **2008**, *36*, D344–50.
- (50) Wishart, D. S.; Knox, C.; Guo, A. C.; Cheng, D.; Shrivastava, S.; Tzur, D.; Gautam, B.; Hassanali, M. DrugBank: a knowledgebase for drugs, drug actions and drug targets. *Nucleic Acids Res.* **2008**, *36*, D901–6.
- (51) Wang, X.; Hao, J.; Yang, Y.; He, K. Natural language adversarial defense through synonym encoding. Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence. 2021; pp 823–833.
- (52) Bjerrum, E. J. SMILES Enumeration as Data Augmentation for Neural Network Modeling of Molecules. *arXiv [cs.LG]* **2017**,
- (53) Gohlke, H.; Mannhold, R.; Kubinyi, H.; Folkers, G. In *Protein-Ligand Interactions*; Gohlke, H., Ed.; Methods and Principles in Medicinal Chemistry; Wiley-VCH Verlag: Weinheim, Germany, 2012.
- (54) Landrum, G. RDKit: A software suite for cheminformatics, computational chemistry, and predictive modeling. http://www.rdkit.org/RDKit_Overview.pdf, 2013; Accessed: 2023-12-13.
- (55) Mukherjee, S.; Ghosh, M.; Basuchowdhuri, P. *Proceedings of the 2022 SIAM International Conference on Data Mining (SDM)*; Proceedings; Society for Industrial and Applied Mathematics, 2022; pp 729–737.
- (56) Chen, L.; Tan, X.; Wang, D.; Zhong, F.; Liu, X.; Yang, T.; Luo, X.; Chen, K.; Jiang, H.; Zheng, M. TransformerCPI: improving compound–protein interaction prediction by sequence-based deep learning with self-attention mechanism and label reversal experiments. *Bioinformatics* **2020**, *36*, 4406–4414.

- (57) Aly Abdelkader, G.; Ngnamsie Njimbouom, S.; Oh, T.-J.; Kim, J.-D. ResBiGAAT: Residual Bi-GRU with attention for protein-ligand binding affinity prediction. *Comput. Biol. Chem.* **2023**, *107*, 107969.
- (58) Li, Q.; Zhang, X.; Wu, L.; Bo, X.; He, S.; Wang, S. PLA-MoRe: A Protein–Ligand Binding Affinity Prediction Model via Comprehensive Molecular Representations. *J. Chem. Inf. Model.* **2022**, *62*, 4380–4390.
- (59) Wang, R.; Fang, X.; Lu, Y.; Wang, S. The PDBbind database: collection of binding affinities for protein-ligand complexes with known three-dimensional structures. *J. Med. Chem.* **2004**, *47*, 2977–2980.
- (60) Chen, S.; Zhang, S.; Fang, X.; Lin, L.; Zhao, H.; Yang, Y. Protein complex structure modeling by cross-modal alignment between cryo-EM maps and protein sequences. *Nat. Commun.* **2024**, *15*, 8808.
- (61) Huang, N.; Shoichet, B. K.; Irwin, J. J. Benchmarking sets for molecular docking. *J. Med. Chem.* **2006**, *49*, 6789–6801.
- (62) Mysinger, M. M.; Carchia, M.; Irwin, J. J.; Shoichet, B. K. Directory of useful decoys, enhanced (DUD-E): better ligands and decoys for better benchmarking. *J. Med. Chem.* **2012**, *55*, 6582–6594.
- (63) Davis, M. I.; Hunt, J. P.; Herrgard, S.; Ciceri, P.; Wodicka, L. M.; Pallares, G.; Hocker, M.; Treiber, D. K.; Zarrinkar, P. P. Comprehensive analysis of kinase inhibitor selectivity. *Nat. Biotechnol.* **2011**, *29*, 1046–1051.
- (64) Tang, J.; Szwajda, A.; Shakyawar, S.; Xu, T.; Hintsanen, P.; Wennerberg, K.; Aitokallio, T. Making sense of large-scale kinase inhibitor bioactivity data sets: a comparative and integrative analysis. *J. Chem. Inf. Model.* **2014**, *54*, 735–743.

- (65) Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T. N.; Weissig, H.; Shindyalov, I. N.; Bourne, P. E. The Protein Data Bank. *Nucleic Acids Res.* **2000**, *28*, 235–242.
- (66) Schomburg, I.; Chang, A.; Schomburg, D. BRENDA, enzyme data and metabolic information. *Nucleic Acids Res.* **2002**, *30*, 47–49.
- (67) Puvanendrapillai, D.; Mitchell, J. B. O. L/D Protein Ligand Database (PLD): additional understanding of the nature and specificity of protein-ligand complexes. *Bioinformatics* **2003**, *19*, 1856–1857.
- (68) Wishart, D. S.; Knox, C.; Guo, A. C.; Shrivastava, S.; Hassanali, M.; Stothard, P.; Chang, Z.; Woolsey, J. DrugBank: a comprehensive resource for in silico drug discovery and exploration. *Nucleic Acids Res.* **2006**, *34*, D668–72.
- (69) Liu, T.; Lin, Y.; Wen, X.; Jorissen, R. N.; Gilson, M. K. BindingDB: a web-accessible database of experimentally determined protein-ligand binding affinities. *Nucleic Acids Res.* **2007**, *35*, D198–201.
- (70) Gaulton, A.; Bellis, L. J.; Bento, A. P.; Chambers, J.; Davies, M.; Hersey, A.; Light, Y.; McGlinchey, S.; Michalovich, D.; Al-Lazikani, B.; Overington, J. P. ChEMBL: a large-scale bioactivity database for drug discovery. *Nucleic Acids Res.* **2012**, *40*, D1100–7.
- (71) Amemiya, T.; Koike, R.; Kidera, A.; Ota, M. PSCDB: a database for protein structural change upon ligand binding. *Nucleic Acids Res.* **2012**, *40*, D554–8.
- (72) Warren, G. L.; Do, T. D.; Kelley, B. P.; Nicholls, A.; Warren, S. D. Essential considerations for using protein-ligand structures in drug discovery. *Drug Discov. Today* **2012**, *17*, 1270–1281.
- (73) Zhu, M.; Song, X.; Chen, P.; Wang, W.; Wang, B. dbHDPLS: A database of human disease-related protein-ligand structures. *Comput. Biol. Chem.* **2019**, *78*, 353–358.

- (74) Wang, C.; Hu, G.; Wang, K.; Brylinski, M.; Xie, L.; Kurgan, L. PDID: database of molecular-level putative protein-drug interactions in the structural human proteome. *Bioinformatics* **2016**, *32*, 579–586.
- (75) Smith, R. D.; Clark, J. J.; Ahmed, A.; Orban, Z. J.; Dunbar, J. B., Jr; Carlson, H. A. Updates to Binding MOAD (Mother of All Databases): Polypharmacology Tools and Their Utility in Drug Repurposing. *J. Mol. Biol.* **2019**, *431*, 2423–2433.
- (76) Gao, M.; Mounbock, A. F. A.; Qaseem, A.; Xu, Q.; Günther, S. CovPDB: a high-resolution coverage of the covalent protein-ligand interactome. *Nucleic Acids Res.* **2022**, *50*, D445–D450.
- (77) Ammar, A.; Cavill, R.; Evelo, C.; Willighagen, E. PSnpBind: a database of mutated binding site protein-ligand complexes constructed using a multithreaded virtual screening workflow. *J. Cheminform.* **2022**, *14*, 8.
- (78) Korlepara, D. B.; Vasavi, C. S.; Jeurkar, S.; Pal, P. K.; Roy, S.; Mehta, S.; Sharma, S.; Kumar, V.; Muvva, C.; Sridharan, B.; Garg, A.; Modee, R.; Bhati, A. P.; Nayar, D.; Priyakumar, U. D. PLAS-5k: Dataset of Protein-Ligand Affinities from Molecular Dynamics for Machine Learning Applications. *Sci Data* **2022**, *9*, 548.
- (79) Lingè, D. et al. PLBD: protein-ligand binding database of thermodynamic and kinetic intrinsic parameters. *Database* **2023**, *2023*.
- (80) Wei, H.; Wang, W.; Peng, Z.; Yang, J. Q-BioLiP: A Comprehensive Resource for Quaternary Structure-based Protein–ligand Interactions. *bioRxiv* **2023**, 2023.06.23.546351.
- (81) Xenarios, I.; Rice, D. W.; Salwinski, L.; Baron, M. K.; Marcotte, E. M.; Eisenberg, D. DIP: the database of interacting proteins. *Nucleic Acids Res.* **2000**, *28*, 289–291.

- (82) Wallach, I.; Lilien, R. The protein-small-molecule database, a non-redundant structural resource for the analysis of protein-ligand binding. *Bioinformatics* **2009**, *25*, 615–620.
- (83) Wang, S.; Lin, H.; Huang, Z.; He, Y.; Deng, X.; Xu, Y.; Pei, J.; Lai, L. CavitySpace: A Database of Potential Ligand Binding Sites in the Human Proteome. *Biomolecules* **2022**, *12*.
- (84) Wang, Y.; You, Z.-H.; Yang, S.; Li, X.; Jiang, T.-H.; Zhou, X. A high efficient biological language model for predicting Protein-Protein interactions. *Cells* **2019**, *8*, 122.
- (85) Liang, W.; KaiYong, Z. Detecting “protein words” through unsupervised word segmentation. *arXiv [cs.CE]* **2014**,
- (86) Abbasi, K.; Razzaghi, P.; Poso, A.; Amanlou, M.; Ghasemi, J. B.; Masoudi-Nejad, A. DeepCDA: deep cross-domain compound–protein affinity prediction through LSTM and convolutional neural networks. *Bioinformatics* **2020**, *36*, 4633–4642.
- (87) Zhou, G.; Gao, Z.; Ding, Q.; Zheng, H.; Xu, H.; Wei, Z.; Zhang, L.; Ke, G. Uni-Mol: A Universal 3D Molecular Representation Learning Framework. *ChemRxiv* **2023**,
- (88) Zhou, D.; Xu, Z.; Li, W.; Xie, X.; Peng, S. MultiDTI: drug–target interaction prediction based on multi-modal representation learning to bridge the gap between new chemical entities and known heterogeneous network. *Bioinformatics* **2021**, *37*, 4485–4492.
- (89) Özçelik, R.; Öztürk, H.; Özgür, A.; Ozkirimli, E. ChemBoost: A chemical language based approach for protein - ligand binding affinity prediction. *Mol. Inform.* **2021**, *40*, e2000212.
- (90) Gaspar, H. A.; Ahmed, M.; Edlich, T.; Fabian, B.; Varszegi, Z.; Segler, M.; Meyers, J.;

- Fiscato, M. Proteochemometric Models Using Multiple Sequence Alignments and a Subword Segmented Masked Language Model. *ChemRxiv* **2021**,
- (91) Arseniev-Koehler, A. Theoretical foundations and limits of word embeddings: What types of meaning can they capture? *Sociol. Methods Res.* **2022**, 004912412211401.
- (92) Lake, B. M.; Murphy, G. L. Word meaning in minds and machines. *Psychol. Rev.* **2023**, *130*, 401–431.
- (93) Winchester, S. A Verb for Our Frantic Times. <https://www.nytimes.com/2011/05/29/opinion/29winchester.html>, 2011; Accessed: 2024-9-15.
- (94) Panapitiya, G.; Girard, M.; Hollas, A.; Sepulveda, J.; Murugesan, V.; Wang, W.; Saldanha, E. Evaluation of deep learning architectures for aqueous solubility prediction. *ACS Omega* **2022**, *7*, 15695–15710.
- (95) Wu, X.; Yu, L. EPSOL: sequence-based protein solubility prediction using multidimensional embedding. *Bioinformatics* **2021**, *37*, 4314–4320.
- (96) Krogh, A. What are artificial neural networks? *Nat. Biotechnol.* **2008**, *26*, 195–197.
- (97) Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; Polosukhin, I. Attention is all you need. *Adv. Neural Inf. Process. Syst.* **2017**, *30*.
- (98) Hochreiter, S.; Schmidhuber, J. Long short-term memory. *Neural Comput.* **1997**, *9*, 1735–1780.
- (99) Sønderby, S. K.; Winther, O. Protein Secondary Structure Prediction with Long Short Term Memory Networks. *arXiv [q-bio.QM]* **2014**,
- (100) Guo, Y.; Li, W.; Wang, B.; Liu, H.; Zhou, D. DeepACLSTM: deep asymmetric convolutional long short-term memory neural models for protein secondary structure prediction. *BMC Bioinformatics* **2019**, *20*, 341.

- (101) Graves, A.; Schmidhuber, J. Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Netw.* **2005**, *18*, 602–610.
- (102) Thafar, M. A.; Alshahrani, M.; Albaradei, S.; Gojobori, T.; Essack, M.; Gao, X. Affinity2Vec: drug-target binding affinity prediction through representation learning, graph mining, and machine learning. *Sci. Rep.* **2022**, *12*, 4751.
- (103) Wei, B.; Zhang, Y.; Gong, X. 519. DeepLPI: A Novel Drug Repurposing Model based on Ligand-Protein Interaction Using Deep Learning. *Open Forum Infect Dis* **2022**, *9*, ofac492.574.
- (104) Yuan, W.; Chen, G.; Chen, C. Y.-C. FusionDTA: attention-based feature polymerizer and knowledge distillation for drug-target binding affinity prediction. *Brief. Bioinform.* **2022**, *23*.
- (105) West-Roberts, J.; Valentin-Alvarado, L.; Mullen, S.; Sachdeva, R.; Smith, J.; Hug, L. A.; Gregoire, D. S.; Liu, W.; Lin, T.-Y.; Husain, G.; Amano, Y.; Ly, L.; Banfield, J. F. Giant genes are rare but implicated in cell wall degradation by predatory bacteria. *bioRxiv* **2023**,
- (106) Vig, J.; Madani, A.; Varshney, L. R.; Xiong, C.; Socher, R.; Rajani, N. F. BERTology Meets Biology: Interpreting Attention in Protein Language Models. *arXiv [cs.CL]* **2020**,
- (107) Koyama, K.; Kamiya, K.; Shimada, K. Cross attention dti: Drug-target interaction prediction with cross attention module in the blind evaluation setup. *BIOKDD2020* **2020**,
- (108) Kurata, H.; Tsukiyama, S. ICAN: Interpretable cross-attention network for identifying drug and target protein interactions. *PLoS One* **2022**, *17*, e0276609.

- (109) Zhao, Q.; Zhao, H.; Zheng, K.; Wang, J. HyperAttentionDTI: improving drug–protein interaction prediction by sequence-based deep learning with attention mechanism. *Bioinformatics* **2021**, *38*, 655–662.
- (110) Wu, H.; Liu, J.; Jiang, T.; Zou, Q.; Qi, S.; Cui, Z.; Tiwari, P.; Ding, Y. AttentionMGT-DTA: A multi-modal drug-target affinity prediction using graph transformer and attention mechanism. *Neural Netw.* **2024**, *169*, 623–636.
- (111) Jiang, M.; Li, Z.; Zhang, S.; Wang, S.; Wang, X.; Yuan, Q.; others Drug–target affinity prediction using graph neural network and contact maps. *scholar.archive.org* **2020**,
- (112) Nguyen, T. M.; Nguyen, T.; Le, T. M.; Tran, T. GEFA: Early Fusion Approach in Drug-Target Affinity Prediction. *IEEE/ACM Trans. Comput. Biol. Bioinform.* **2022**, *19*, 718–728.
- (113) Yu, J.; Li, Z.; Chen, G.; Kong, X.; Hu, J.; Wang, D.; Cao, D.; Li, Y.; Huo, R.; Wang, G.; Liu, X.; Jiang, H.; Li, X.; Luo, X.; Zheng, M. Computing the relative binding affinity of ligands based on a pairwise binding comparison network. *Nature Computational Science* **2023**, *3*, 860–872.
- (114) Knutson, C.; Bontha, M.; Bilbrey, J. A.; Kumar, N. Decoding the protein–ligand interactions using parallel graph neural networks. *Sci. Rep.* **2022**, *12*, 1–14.
- (115) Kyro, G. W.; Brent, R. I.; Batista, V. S. HAC-Net: A Hybrid Attention-Based Convolutional Neural Network for Highly Accurate Protein–Ligand Binding Affinity Prediction. *J. Chem. Inf. Model.* **2023**, *63*, 1947–1960.
- (116) Yousefi, N.; Yazdani-Jahromi, M.; Tayebi, A.; Kolanthai, E.; Neal, C. J.; Banerjee, T.; Gosai, A.; Balasubramanian, G.; Seal, S.; Ozmen Garibay, O. BindingSite-AugmentedDTA: enabling a next-generation pipeline for interpretable prediction models in drug repurposing. *Brief. Bioinform.* **2023**, *24*.

- (117) Yazdani-Jahromi, M.; Yousefi, N.; Tayebi, A.; Kolanthai, E.; Neal, C. J.; Seal, S.; Garibay, O. O. AttentionSiteDTI: an interpretable graph-based model for drug-target interaction prediction using NLP sentence-level relation classification. *Brief. Bioinform.* **2022**, *23*.
- (118) Lim, J.; Ryu, S.; Park, K.; Choe, Y. J.; Ham, J.; Kim, W. Y. Predicting Drug-Target Interaction Using a Novel Graph Neural Network with 3D Structure-Embedded Graph Representation. *J. Chem. Inf. Model.* **2019**, *59*, 3981–3988.
- (119) Lin, Z.; Akin, H.; Rao, R.; Hie, B.; Zhu, Z.; Lu, W.; Smetanin, N.; Verkuil, R.; Kabeli, O.; Shmueli, Y.; Dos Santos Costa, A.; Fazel-Zarandi, M.; Sercu, T.; Candido, S.; Rives, A. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science* **2023**, *379*, 1123–1130.
- (120) Zhang, S.; Fan, R.; Liu, Y.; Chen, S.; Liu, Q.; Zeng, W. Applications of transformer-based language models in bioinformatics: a survey. *Bioinform. Adv.* **2023**, *3*, vbad001.
- (121) Bahdanau, D.; Cho, K.; Bengio, Y. Neural Machine Translation by Jointly Learning to Align and Translate. *arXiv [cs.CL]* **2014**,
- (122) Zeyer, A.; Bahar, P.; Irie, K.; Schlüter, R.; Ney, H. A Comparison of Transformer and LSTM Encoder Decoder Models for ASR. 2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU). 2019; pp 8–15.
- (123) Irie, K.; Zeyer, A.; Schlüter, R.; Ney, H. Language Modeling with Deep Transformers. *arXiv [cs.CL]* **2019**,
- (124) Zouitni, C.; Sabri, M. A.; Aarab, A. A Comparison Between LSTM and Transformers for Image Captioning. *Digital Technologies and Applications*. 2023; pp 492–500.
- (125) Parisotto, E.; Song, F.; Rae, J.; Pascanu, R.; Gulcehre, C.; Jayakumar, S.; Jaderberg, M.; Kaufman, R. L.; Clark, A.; Noury, S.; Botvinick, M.; Heess, N.; Hadsell, R.

- Stabilizing Transformers for Reinforcement Learning. Proceedings of the 37th International Conference on Machine Learning. 2020; pp 7487–7498.
- (126) Bilokon, P.; Qiu, Y. Transformers versus LSTMs for electronic trading. *arXiv [q-fin.TR]* **2023**,
- (127) Merity, S. Single Headed Attention RNN: Stop Thinking With Your Head. *arXiv [cs.CL]* **2019**,
- (128) Ezen-Can, A. A Comparison of LSTM and BERT for Small Corpus. *arXiv [cs.CL]* **2020**,
- (129) Unsal, S.; Atas, H.; Albayrak, M.; Turhan, K.; Acar, A. C.; Doğan, T. Learning functional properties of proteins with language models. *Nature Machine Intelligence* **2022**, *4*, 227–245.
- (130) Brandes, N.; Ofer, D.; Peleg, Y.; Rappoport, N.; Linial, M. ProteinBERT: a universal deep-learning model of protein sequence and function. *Bioinformatics* **2022**, *38*, 2102–2110.
- (131) Luo, S.; Chen, T.; Xu, Y.; Zheng, S.; Liu, T.-Y.; Wang, L.; He, D. One Transformer Can Understand Both 2D & 3D Molecular Data. *arXiv [cs.LG]* **2022**,
- (132) Devlin, J.; Chang, M.-W.; Lee, K.; Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv [cs.CL]* **2018**,
- (133) Clark, K.; Luong, M.-T.; Le, Q. V.; Manning, C. D. ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators. *arXiv [cs.CL]* **2020**,
- (134) Wang, J.; Wen, N.; Wang, C.; Zhao, L.; Cheng, L. ELECTRA-DTA: a new compound-protein binding affinity prediction model based on the contextualized sequence encoding. *J. Cheminform.* **2022**, *14*, 14.

- (135) Shin, B.; Park, S.; Kang, K.; Ho, J. C. Self-Attention Based Molecule Representation for Predicting Drug-Target Interaction. Proceedings of the 4th Machine Learning for Healthcare Conference. 2019; pp 230–248.
- (136) Huang, K.; Xiao, C.; Glass, L. M.; Sun, J. MolTrans: Molecular Interaction Transformer for drug–target interaction prediction. *Bioinformatics* **2020**, *37*, 830–836.
- (137) Shen, L.; Feng, H.; Qiu, Y.; Wei, G.-W. SVSBI: sequence-based virtual screening of biomolecular interactions. *Commun Biol* **2023**, *6*, 536.
- (138) Wang, J.; Hu, J.; Sun, H.; Xu, M.; Yu, Y.; Liu, Y.; Cheng, L. MGPLI: exploring multigranular representations for protein–ligand interaction prediction. *Bioinformatics* **2022**, *38*, 4859–4867.
- (139) Qian, Y.; Wu, J.; Zhang, Q. CAT-CPI: Combining CNN and transformer to learn compound image features for predicting compound-protein interactions. *Front Mol Biosci* **2022**, *9*, 963912.
- (140) Nguyen, T.; Le, H.; Quinn, T. P.; Nguyen, T.; Le, T. D.; Venkatesh, S. GraphDTA: predicting drug–target binding affinity with graph neural networks. *Bioinformatics* **2020**, *37*, 1140–1147.
- (141) Jin, Z.; Wu, T.; Chen, T.; Pan, D.; Wang, X.; Xie, J.; Quan, L.; Lyu, Q. CAPLA: improved prediction of protein–ligand binding affinity by a deep learning approach based on a cross-attention mechanism. *Bioinformatics* **2023**, *39*, btad049.
- (142) Nam, H.; Ha, J.-W.; Kim, J. Dual attention networks for multimodal reasoning and matching. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). 2017; pp 299–307.
- (143) Wang, X.; Liu, D.; Zhu, J.; Rodriguez-Paton, A.; Song, T. CSConv2d: A 2-D Struc-

- tural Convolution Neural Network with a Channel and Spatial Attention Mechanism for Protein-Ligand Binding Affinity Prediction. *Biomolecules* **2021**, *11*.
- (144) Anteghini, M.; Santos, V. A. M. D.; Saccenti, E. PortPred: Exploiting deep learning embeddings of amino acid sequences for the identification of transporter proteins and their substrates. *J. Cell. Biochem.* **2023**,
- (145) Huang, K.; Fu, T.; Glass, L. M.; Zitnik, M.; Xiao, C.; Sun, J. DeepPurpose: a deep learning library for drug–target interaction prediction. *Bioinformatics* **2020**, *36*, 5545–5547.
- (146) Zhao, L.; Wang, J.; Pang, L.; Liu, Y.; Zhang, J. GANsDTA: Predicting Drug-Target Binding Affinity Using GANs. *Front. Genet.* **2019**, *10*, 1243.
- (147) Hu, F.; Jiang, J.; Wang, D.; Zhu, M.; Yin, P. Multi-PLI: interpretable multi-task deep learning model for unifying protein–ligand interaction datasets. *J. Cheminform.* **2021**, *13*, 30.
- (148) Zheng, S.; Li, Y.; Chen, S.; Xu, J.; Yang, Y. Predicting Drug Protein Interaction using Quasi-Visual Question Answering System. *bioRxiv* **2019**, 588178.
- (149) Tsubaki, M.; Tomii, K.; Sese, J. Compound–protein interaction prediction with end-to-end learning of neural networks for graphs and sequences. *Bioinformatics* **2018**, *35*, 309–318.
- (150) Karimi, M.; Wu, D.; Wang, Z.; Shen, Y. DeepAffinity: interpretable deep learning of compound–protein affinity through unified recurrent and convolutional neural networks. *Bioinformatics* **2019**, *35*, 3329–3338.
- (151) Li, S.; Wan, F.; Shu, H.; Jiang, T.; Zhao, D.; Zeng, J. MONN: a multi-objective neural network for predicting compound–protein interactions and affinities. *Cell Systems* **2020**, *10*, 308–322.

- (152) Zhao, M.; Yuan, M.; Yang, Y.; Xu, S. X. CPGL: Prediction of Compound-Protein Interaction by Integrating Graph Attention Network With Long Short-Term Memory Neural Network. *IEEE/ACM Trans. Comput. Biol. Bioinform.* **2023**, *20*, 1935–1942.
- (153) Yu, L.; Qiu, W.; Lin, W.; Cheng, X.; Xiao, X.; Dai, J. HGDTI: predicting drug–target interaction by using information aggregation based on heterogeneous graph neural network. *BMC Bioinformatics* **2022**, *23*, 126.
- (154) Lee, I.; Nam, H. Sequence-based prediction of protein binding regions and drug-target interactions. *J. Cheminform.* **2022**, *14*, 5.
- (155) Gönen, M.; Heller, G. Concordance probability and discriminatory power in proportional hazards regression. *Biometrika* **2005**, *92*, 965–970.
- (156) M., C. *Pattern Recognition and Machine Learning*, 1st ed.; Information Science and Statistics; Springer: New York, NY, 2006.
- (157) Deller, M. C.; Rupp, B. Models of protein-ligand crystal structures: trust, but verify. *J. Comput. Aided Mol. Des.* **2015**, *29*, 817–836.
- (158) Nasteski, V. An overview of the supervised machine learning methods. *Horizons* **2017**, *4*, 51–62.
- (159) Kozlov, M. So you got a null result. Will anyone publish it? *Nature* **2024**, *631*, 728–730.
- (160) Edfeldt, K. et al. A data science roadmap for open science organizations engaged in early-stage drug discovery. *Nat. Commun.* **2024**, *15*, 5640.
- (161) Mlinarić, A.; Horvat, M.; Šupak Smolčić, V. Dealing with the positive publication bias: Why you should really publish your negative results. *Biochem. Med.* **2017**, *27*, 030201.

- (162) Albalade, A.; Minker, W. *Semi-supervised and unsupervised machine learning: Novel strategies*; Wiley-ISTE, 2013.
- (163) Sajadi, S. Z.; Zare Chahooki, M. A.; Gharaghani, S.; Abbasi, K. AutoDTI++: deep unsupervised learning for DTI prediction by autoencoders. *BMC Bioinformatics* **2021**, *22*, 204.
- (164) Najm, M.; Azencott, C.-A.; Playe, B.; Stoven, V. Drug Target Identification with Machine Learning: How to Choose Negative Examples. *Int. J. Mol. Sci.* **2021**, *22*.
- (165) Yang, J.; Shen, C.; Huang, N. Predicting or Pretending: Artificial Intelligence for Protein-Ligand Interactions Lack of Sufficiently Large and Unbiased Datasets. *Front. Pharmacol.* **2020**, *11*, 69.
- (166) Sieg, J.; Flachsenberg, F.; Rarey, M. In Need of Bias Control: Evaluating Chemical Data for Machine Learning in Structure-Based Virtual Screening. *J. Chem. Inf. Model.* **2019**, *59*, 947–961.
- (167) Volkov, M.; Turk, J.-A.; Drizard, N.; Martin, N.; Hoffmann, B.; Gaston-Mathé, Y.; Rognan, D. On the Frustration to Predict Binding Affinities from Protein–Ligand Structures with Deep Neural Networks. *J. Med. Chem.* **2022**, *65*, 7946–7958.
- (168) Takahashi, M.; Maraboeuf, F.; Nordén, B. Locations of functional domains in the RecA protein. Overlap of domains and regulation of activities. *Eur. J. Biochem.* **1996**, *242*, 20–28.
- (169) Lee, I.; Nam, H. Infusing Linguistic Knowledge of SMILES into Chemical Language Models. *arXiv [q-bio.QM]* **2022**,
- (170) Skinnider, M. A. Invalid SMILES are beneficial rather than detrimental to chemical language models. *Nature Machine Intelligence* **2024**, 1–12.

- (171) O’Boyle, N.; Dalke, A. DeepSMILES: An Adaptation of SMILES for Use in Machine-Learning of Chemical Structures. *ChemRxiv* **2018**,
- (172) Krenn, M.; Häse, F.; Nigam, A.; Friederich, P.; Aspuru-Guzik, A. Self-referencing embedded strings (SELFIES): A 100% robust molecular string representation. *Mach. Learn.: Sci. Technol.* **2020**, *1*, 045024.
- (173) Bibal, A.; Cardon, R.; Alfter, D.; Wilkens, R.; Wang, X.; François, T.; Watrin, P. Is Attention Explanation? An Introduction to the Debate. Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Dublin, Ireland, 2022; pp 3889–3900.
- (174) Lundberg, S. M.; Lee, S.-I. A unified approach to interpreting model predictions. *Neural Inf Process Syst* **2017**, *30*, 4765–4774.
- (175) Gu, Y.; Zhang, X.; Xu, A.; Chen, W.; Liu, K.; Wu, L.; Mo, S.; Hu, Y.; Liu, M.; Luo, Q. Protein-ligand binding affinity prediction with edge awareness and supervised attention. *iScience* **2023**, *26*, 105892.
- (176) Rodis, N.; Sardianos, C.; Papadopoulos, G. T.; Radoglou-Grammatikis, P.; Sarigiannidis, P.; Varlamis, I. Multimodal Explainable Artificial Intelligence: A Comprehensive Review of Methodological Advances and Future Research Directions. *arXiv [cs.AI]* **2023**,
- (177) Gilpin, L. H.; Bau, D.; Yuan, B. Z.; Bajwa, A.; Specter, M.; Kagal, L. Explaining explanations: An overview of interpretability of machine learning. 2018 IEEE 5th International Conference on Data Science and Advanced Analytics (DSAA). 2018; pp 80–89.
- (178) Popel, M.; Bojar, O. Training tips for the Transformer model. *arXiv [cs.CL]* **2018**,

- (179) Eguida, M.; Rognan, D. A Computer Vision Approach to Align and Compare Protein Cavities: Application to Fragment-Based Drug Design. *J. Med. Chem.* **2020**, *63*, 7127–7142.
- (180) Öztürk, H.; Özgür, A.; Ozkirimli, E. DeepDTA: deep drug-target binding affinity prediction. *Bioinformatics* **2018**, *34*, i821–i829.
- (181) Park, Y.; Metzger, B. P. H.; Thornton, J. W. The simplicity of protein sequence-function relationships. *Nat. Commun.* **2024**, *15*, 7953.
- (182) Evans, R. et al. Protein complex prediction with AlphaFold-Multimer. *bioRxiv* **2021**,
- (183) Stahl, K.; Warneke, R.; Demann, L.; Bremenkamp, R.; Hormes, B.; Brock, O.; Stülke, J.; Rappsilber, J. Modelling protein complexes with crosslinking mass spectrometry and deep learning. *Nat. Commun.* **2024**, *15*, 7866.
- (184) Kryshafaovych, A.; Schwede, T.; Topf, M.; Fidelis, K.; Moulton, J. Critical assessment of methods of protein structure prediction (CASP)—Round XIII. *Proteins* **2019**, *87*, 1011–1020.
- (185) Janin, J.; Henrick, K.; Moulton, J.; Eyck, L. T.; Sternberg, M. J. E.; Vajda, S.; Vakser, I.; Wodak, S. J. CAPRI: A Critical Assessment of PRedicted Interactions. *Proteins* **2003**, *52*, 2–9.
- (186) Lensink, M. F.; Nadzirin, N.; Velankar, S.; Wodak, S. J. Modeling protein-protein, protein-peptide, and protein-oligosaccharide complexes: CAPRI 7th edition. *Proteins* **2020**, *88*, 916–938.
- (187) Senior, A. W. et al. Improved protein structure prediction using potentials from deep learning. *Nature* **2020**, *577*, 706–710.
- (188) Baek, M. et al. Accurate prediction of protein structures and interactions using a three-track neural network. *Science* **2021**, *373*, 871–876.

- (189) Kramer, M. A. Nonlinear principal component analysis using autoassociative neural networks. *AIChE J.* **1991**, *37*, 233–243.
- (190) Rogers, D.; Hahn, M. Extended-connectivity fingerprints. *J. Chem. Inf. Model.* **2010**, *50*, 742–754.
- (191) Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Liò, P.; Bengio, Y. Graph Attention Networks. *arXiv [stat.ML]* **2017**,
- (192) Kipf, T. N.; Welling, M. Semi-Supervised Classification with Graph Convolutional Networks. *arXiv [cs.LG]* **2016**,
- (193) Xu, Z.; Wang, S.; Zhu, F.; Huang, J. Seq2seq Fingerprint: An Unsupervised Deep Molecular Embedding for Drug Discovery. Proceedings of the 8th ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics. New York, NY, USA, 2017; pp 285–294.
- (194) Cho, K.; van Merriënboer, B.; Gulcehre, C.; Bahdanau, D.; Bougares, F.; Schwenk, H.; Bengio, Y. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. *arXiv [cs.CL]* **2014**,
- (195) Gilmer, J.; Schoenholz, S.; Riley, P. F.; Vinyals, O.; Dahl, G. E. Neural Message Passing for Quantum Chemistry. *ICML* **2017**, 1263–1272.
- (196) Asgari, E.; Mofrad, M. R. K. Continuous Distributed Representation of Biological Sequences for Deep Proteomics and Genomics. *PLoS One* **2015**, *10*, e0141287.
- (197) He, K.; Zhang, X.; Ren, S.; Sun, J. Deep residual learning for image recognition. *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.* **2015**, 770–778.
- (198) Öztürk, H.; Ozkirimli, E.; Özgür, A. A novel methodology on distributed representations of proteins using their interacting ligands. *Bioinformatics* **2018**, *34*, i295–i303.

- (199) Hu, J.; Shen, L.; Sun, G. Squeeze-and-excitation networks. Proceedings of the IEEE conference on computer vision and pattern recognition. 2018; pp 7132–7141.

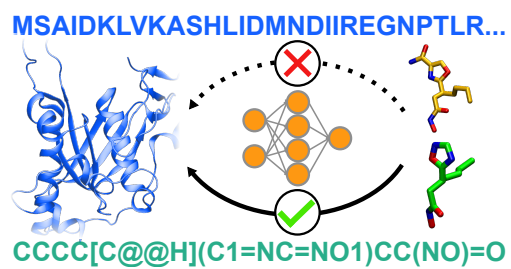


Table of Contents Graphic