# Advances in Protein Representation Learning: Methods, Applications, and Future Directions

Viet Thanh Duy Nguyen and Truong-Son Hy*

*Department of Computer Science, The University of Alabama at Birmingham, Birmingham, AL 35294, United States*

E-mail: thy@uab.edu

**Abstract**

Proteins are complex biomolecules that play a central role in various biological processes, making them critical targets for breakthroughs in molecular biology, medical research, and drug discovery. Deciphering their intricate, hierarchical structures, and diverse functions is essential for advancing our understanding of life at the molecular level. Protein Representation Learning (PRL) has emerged as a transformative approach, enabling the extraction of meaningful computational representations from protein data to address these challenges. In this paper, we provide a comprehensive review of PRL research, categorizing methodologies into five key areas: feature-based, sequence-based, structure-based, multimodal, and complex-based approaches. To support researchers in this rapidly evolving field, we introduce widely used databases for protein sequences, structures, and functions, which serve as essential resources for model development and evaluation. We also explore the diverse applications of these approaches in multiple domains, demonstrating their broad impact. Finally, we discuss pressing technical challenges and outline future directions to advance PRL, offering insights to inspire continued innovation in this foundational field.

# 1. Introduction

Proteins are fundamental biomolecules that are responsible for a wide range of biological processes, including enzymatic catalysis, structural support, signal transduction, and molecular recognition. Their function is dictated by a hierarchical structure comprising four levels: primary (amino acid sequence), secondary (local folding patterns), tertiary (three-dimensional conformation), and quaternary (multi-subunit assembly), as shown in Fig. 1. Understanding these structures is essential for applications in molecular biology, drug discovery, and biotechnology. Understanding protein structure and function is essential to advance molecular biology, drug discovery, and biotechnology. However, the vast diversity of proteins and the complexity of their interactions pose significant challenges to traditional computational and experimental approaches. Protein Representation Learning (PRL) has emerged as a transformative tool for addressing these challenges by encoding protein sequences, structures, and functions into compact, meaningful representations that can be used for various predictive and generative tasks.

Over the past decade, PRL has witnessed remarkable advancements, fueled by innovations in machine learning and the growing availability of large-scale protein datasets. From feature-based encoding methods to advanced multimodal frameworks that integrate sequence, structure, and functional data, PRL has enabled significant breakthroughs in protein property prediction, structure modeling, and design. These developments have not only advanced fundamental protein science but also opened new avenues for applications in therapeutic development, enzyme engineering, and biomolecular design.

Although numerous surveys have been conducted in the field of PRL, most have been developed from the perspective of biological applications, model architectures, or pretext tasks.[1–3] Although these reviews provide valuable information, they often focus on specialized aspects of PRL, making them less accessible and applicable to researchers outside of AI and computer science. In contrast, our review categorizes PRL research based on the specific modalities utilized, such as sequence, structure, and function, offering a structured framework for understanding how these data types are applied in different approaches. As illustrated in Fig. 2, this organization
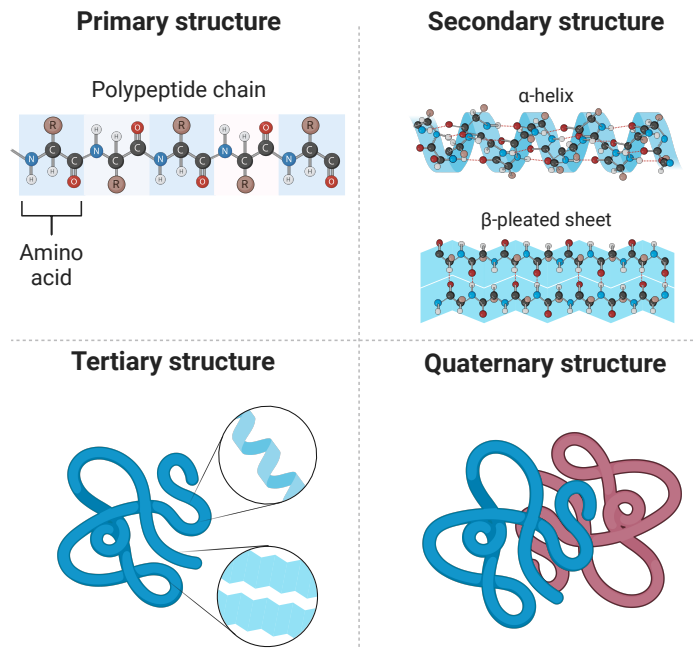
Figure 1: The four levels of protein structure, organized by increasing complexity within the polypeptide chain. Primary structure refers to the specific sequence of amino acids. Secondary structure involves local folding patterns, such as $\alpha$-helices and $\beta$-sheets. Tertiary structure represents the overall three-dimensional conformation of a single polypeptide chain. Quaternary structure describes the assembly and interactions of multiple polypeptide chains within a protein complex.

provides a clear pathway for researchers in various fields to identify PRL methods that align with their data types and research needs. By emphasizing the role of each modality and its relevance in diverse domains, this review bridges the gap between computational advancements and real-world applications, serving as a valuable resource for interdisciplinary researchers seeking to leverage PRL effectively.

In general, the structure of this paper is organized as follows: it begins with Section 1. Introduction, which provides context on PRL and its significance in understanding the structure and function of the protein. Section 2. Feature-Based Approaches explores early methods that encode proteins on the basis of their physical properties. Section 3. Sequence-Based Approaches discusses both non-aligned and aligned sequence methods for capturing protein sequence information. Section 4. Structure-Based Approaches delves into residue-level, atomic-level, protein surface representations, along with symmetry-preserving and equivariant representation learning,

emphasizing the critical role of structural data in PRL. Section 5. Multimodal-Based Approaches examines frameworks that integrate sequence, structure, and functional data to generate enriched representations. Section 6. Complex-Based Approaches covers protein-ligand and protein-protein complex representations, focusing on modeling molecular interactions. To facilitate comparison, we provide a structured overview of the strengths and limitations of each approach in Table 1, offering insights to guide the selection of appropriate PRL strategies based on specific biological and computational requirements. Section 7. Databases for Protein Representation Learning introduces key resources for sequences, structures, and functions that support PRL. Section 8. Applications highlights practical implementations of PRL, including protein property prediction, protein structure prediction, protein design and optimization (e.g., ligand-binding proteins, enzymes, and antibodies), and structure-based drug design. Section 9. Future Directions and Open Challenges discusses key challenges such as expanding PRL to DNA / RNA representation learning, improving the scalability and generalization of the model, and improving the explainability, with the aim of inspiring future research and advances in the field. Finally, Section 10. Conclusion provides a summary of the key insights from this review and reflects on the broader impact of PRL, emphasizing its potential for advancing protein science and biomedical applications.

## 2. Feature-Based Approaches

Early computational representations of proteins relied on encoding amino acid sequences using physicochemical and biochemical properties.[4,5] This approach assigns numerical attributes—such as charge, hydrophobicity, and molecular size—to amino acids, often sourced from standardized databases like AAIndex.[6] These numerical feature vectors enable traditional machine learning models to analyze protein properties based on their physicochemical traits.

Despite their utility, feature-based methods face challenges:[7] (i) Selecting relevant properties is non-trivial, as molecular behavior is highly context-dependent and often poorly understood; and (ii) Handcrafted feature engineering may oversimplify protein characteristics, limiting its ability to
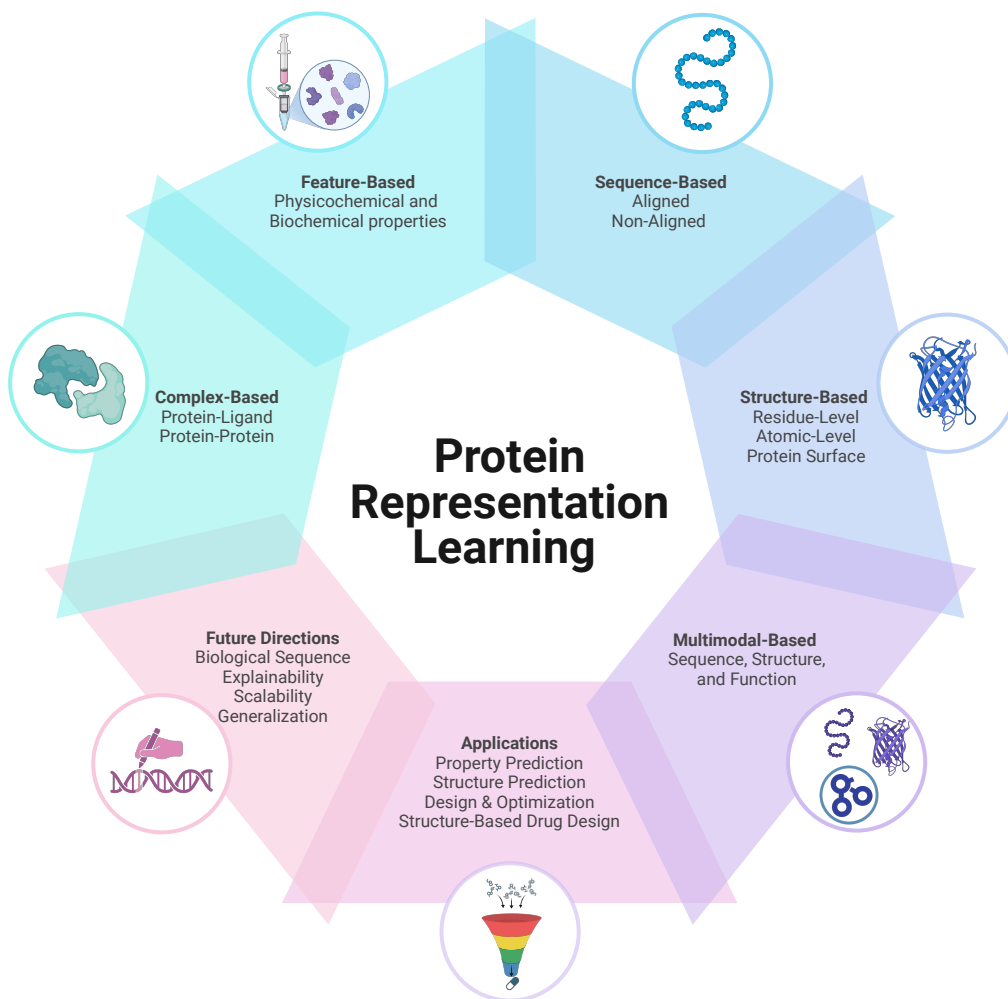
Figure 2: An overview of the key components and themes discussed in this review. The figure highlights the interconnections and relationships between the main topics, providing a comprehensive visual summary of the review's scope.

capture complex sequence-function relationships. As a result, these representations often struggle to generalize across diverse protein-related tasks.

To address these limitations, modern PRL approaches have shifted toward data-driven methods that automatically extract informative representations from raw protein sequences or structures. By leveraging machine learning, these techniques eliminate the need for manual feature selection, improving scalability and predictive performance in applications such as drug discovery, functional annotation, and protein engineering.

# 3. Sequence-Based Approaches

Protein sequences, composed of linear chains of amino acids, serve as the foundation for under-standing protein behavior. These sequences can be viewed as a unique "biological language", where amino acids act as the "words" arranged in a linear "sentence". This analogy to natural (tex-tual) language has driven the adoption of methodologies from natural language processing (NLP) for PRL, enabling models to capture the implicit "grammar" encoded within protein sequences. Sequence-based approaches focus on extracting meaningful representations directly from these sequences and can be broadly categorized into two main methods: aligned approaches, which utilize evolutionary information from Multiple Sequence Alignments (MSAs), and non-aligned approaches, which process individual sequences independently.

## 3.1. Aligned Sequence Approaches

Aligned sequence methods leverage evolutionary relationships by analyzing Multiple Sequence Alignments (MSAs), which align homologous sequences to highlight conserved regions critical for protein function. The long-standing practice in computational biology is to infer structural and functional constraints from evolutionarily related sequences, as conserved regions often correspond to functionally or structurally important sites. Consequently, several protein representation models have been developed to capture co-evolutionary information by taking MSAs as input. These approaches have been particularly effective for protein folding and structure prediction, where evolutionary information plays a critical role.[8–11]

One of the most influential works leveraging MSAs is AlphaFold,[8] which marked a significant breakthrough in addressing the protein folding problem. By delivering unprecedented accuracy in predicting protein structures, AlphaFold fundamentally advanced the field and established a foundation for numerous subsequent research endeavors, inspiring innovations in protein structure prediction and related applications.

However, MSA-based methods have limitations. First, they rely on the availability of homol-

ogous sequences, which may not exist for every protein.[12] Second, some proteins, such as those from orphan families or de novo-designed proteins, lack sufficient evolutionary data, limiting the applicability of these models.[12] Additionally, the computational cost of constructing MSAs for large protein families can be prohibitive, leading researchers to explore non-aligned sequence approaches.[13]

## 3.2. Non-Aligned Sequence Approaches

Non-aligned methods analyze protein sequences independently, without incorporating evolutionary relationships. These approaches are particularly useful for proteins without sufficient homologous sequences and provide an alternative when MSAs are computationally infeasible.

Early non-aligned approaches relied on handcrafted sequence features, such as amino acid composition, k-mer frequency, physicochemical properties, and evolutionary scoring matrices. These features were fed into classical machine learning models, including Support Vector Machines (SVMs) and Random Forests, for protein classification and function prediction.[14,15]

The advent of deep learning led to the adoption of Variational Autoencoders (VAEs), which learn low-dimensional latent representations of protein sequences, enabling efficient exploration of sequence space.[16,17] Other works employed Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks to capture sequential dependencies within protein sequences, improving sequence modeling tasks.[18] Additionally, unsupervised Doc2Vec embedding models have been used to learn fixed-length representations of protein sequences by treating them as "documents" composed of "words".[7]

More recently, the field has progressed toward transformer-based protein language models, driven by advancements in large language models (LLMs) from natural language processing. These Protein Language Models (PLMs), such as Evolutionary Scale Modeling (ESM) series,[13,19,20] ProtTrans,[21] ProteinBERT,[22] and ProLLaMA,[23] have been pre-trained on large-scale protein sequence datasets. Unlike earlier models, PLMs capture contextual relationships between amino acids, enabling the generation of rich embeddings that encode biologically meaningful informa-

tion. One major advantage of PLMs is that they generalize beyond known protein families, making them effective for de novo protein design and function prediction.

Beyond general-purpose PLMs, specialized models have been developed for specific functional protein families, such as antibody-specific models[24–27] and enzyme-focused models[28] trained to capture catalytic properties and functional specificity. These specialized PLMs incorporate domain-specific inductive biases, making them more effective for tasks such as antibody affinity maturation, enzyme engineering, and therapeutic protein design.

# 4. Structure-Based Approaches

Structure-based approaches capture the three-dimensional organization of proteins which is crucial for understanding their structural and functional properties. These approaches leverage structural data at various levels of granularity, as illustrated in Fig. 3, including atoms, residues, and surfaces, to provide richer insights into protein stability, function, and interactions. Additionally, they integrate symmetry principles to enhance learning efficiency and improve representation quality.

## 4.1. Residue-Level Representation

Residue-level representation models focus on capturing structural information at the amino acid level. While sequence-based representations are limited to linear order, residue-level approaches consider the spatial arrangement of residues in the folded protein. This is particularly important because residues that are distant in the primary sequence may be brought into close proximity in the three-dimensional structure due to protein folding. By explicitly modeling these spatial relationships, residue-level representations enable a more comprehensive understanding of protein stability, interactions, and functional properties.

Several notable studies have advanced residue-level representations by leveraging graph-based approaches.[29–38] In this representation, proteins are modeled as graphs, where nodes correspond to residues, typically represented by the alpha carbon ($C\alpha$) of each amino acid along with its 3D
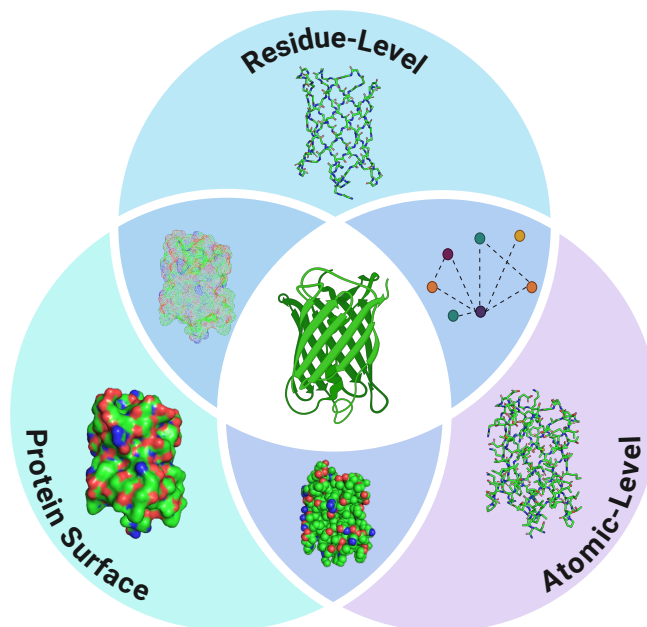
8

Figure 3: Illustration of different structural representations of a protein at the tertiary level, along with their typical computational representations. Residue-Level Representation models the protein backbone, typically using alpha carbon (C$\alpha$) atoms to capture the overall fold and residue connectivity, with a corresponding graph-based representation. Atomic-Level Representation considers all individual atoms within the protein structure, including backbone and side-chain atoms, commonly represented as a point cloud or an all-atom graph that captures atomic interactions. Protein Surface Representation focuses on the solvent-accessible surface, highlighting geometric and physicochemical properties that influence binding and molecular recognition, often modeled using a surface mesh or a point cloud representation that encodes local curvature and electrostatic potential.

coordinates. Node features frequently include the amino acid type, physicochemical properties (e.g., charge, hydrophobicity, and polarity), and secondary structure annotations, providing rich information about each residue's role and environment. The edges of the residue-level graphs vary to reflect different relationships between residues:[36]

- **Sequential edges** connect residues that are adjacent in the primary sequence, capturing the linear backbone connectivity of the protein chain.

- **Radius edges** link residues within a predefined distance cut-off in 3D space, emphasizing local spatial neighborhoods critical for structural integrity.

- **K-nearest neighbor edges** connect each residue to its closest spatial neighbors, independent of their sequence order, effectively capturing geometric relationships across the protein.

9

Additionally, edge features may encode interaction types, such as hydrogen bonds, disulfide bonds, or hydrophobic interactions, as well as geometric properties like bond angles and dihedral angles, further enhancing the model's ability to capture structural and functional relationships.

## 4.2. Atomic-Level Representation

Atomic-level representation models capture the fine-grained details of protein structures by encoding the positions and interactions of individual atoms, offering a more detailed view than residue-level approaches. While residue-level representations typically approximate each amino acid with a single point - often the alpha carbon ($C\alpha$) - atomic-level models consider all atoms, including side chains and backbone atoms, preserving essential structural and chemical properties. Side-chain atoms play a fundamental role in molecular recognition, enzymatic activity, and ligand binding, where interactions occur at the atomic scale. By explicitly modeling all atoms, atomic-level representations enhance the ability to predict fine-grained molecular interactions.

Several notable studies have advanced atomic-level representations by leveraging graph-based[34,39,40] and point cloud-based approaches.[41,42] In graph-based models, proteins are represented as atomic-level graphs, where nodes correspond to individual atoms, each characterized by features such as atomic type, physicochemical properties, and the types of amino acids that the atom belongs to. Edges capture interactions between atoms, commonly defined based on chemical bonds, non-covalent interactions, or k-nearest neighbors (KNN), which connect each atom to its closest spatial neighbors in 3D space. The KNN-based approach, which is the most commonly used edge construction method in atomic-level graphs, allows models to incorporate both bonded and non-bonded atomic interactions, improving their ability to capture structural and functional relationships. In contrast, point cloud-based representations treat atoms as independent points in three-dimensional space, discarding explicit connectivity but preserving raw spatial distributions. This approach enables models to learn atomic-level geometric features without requiring predefined graph structures, making them particularly effective for capturing fine-grained structural details and local atomic interactions.

## 4.3. Protein Surface Representation

A protein's molecular surface is a compact, smooth boundary formed by the outermost atoms, exhibiting distinct chemical and geometric patterns. As the primary interface for molecular interactions, the protein surface plays a crucial role in molecular recognition, ligand binding, and protein-protein interactions. Protein surface representation focuses on modeling this exterior, where these interactions predominantly occur, providing a functionally relevant perspective that complements residue-level and atomic-level representations.

Several notable studies have advanced protein surface representations by leveraging 3D mesh-based approaches[43–47] and point cloud-based approaches.[48–51] In mesh-based models, the protein surface is represented as a triangulated mesh, where nodes correspond to discrete surface points, and edges define geometric relationships between neighboring points, preserving the overall topology of the surface. These models effectively capture geometric features such as curvature, solvent accessibility, and surface normals, as well as chemical properties like electrostatic potential and hydrophobicity. Due to their structured nature, mesh-based representations excel in tasks that require fine-grained surface characterization, such as functional site identification, enzyme active site modeling, and predicting molecular docking orientations. However, mesh-based approaches often rely on precomputed geometric and chemical features, making them computationally demanding and less adaptable for large-scale datasets.

In contrast, point cloud-based approaches offer a more computationally efficient alternative by avoiding the need for precomputed features. These models operate directly on the raw set of atoms composing the protein, generating a point cloud representation of the surface on the fly. Unlike mesh-based methods, point cloud approaches learn task-specific geometric and chemical features dynamically and apply convolutional operators that approximate geodesic coordinates in the tangent space. This eliminates the need for an explicit surface mesh and significantly reduces memory usage and computational overhead, making point cloud-based representations well-suited for large-scale protein-ligand binding predictions, protein-protein interface modeling, and high-throughput interaction analysis.

## 4.4. Symmetry-Preservation and Equivariance in Structural Representation

Incorporating three-dimensional (3D) structural information into PRL presents the challenge of ensuring models correctly handle spatial transformations such as rotations and translations. Models relying on raw 3D coordinates risk developing spurious dependencies on absolute positioning, leading to inconsistent predictions for identical structures in different orientations. This limitation is particularly problematic in tasks like protein folding, molecular docking, and binding site prediction, where relative spatial geometry governs molecular behavior. Symmetry-preserving and equivariant models addresses this issue by explicitly incorporating the inherent rotational and translational symmetries of protein structures, improving model performance, generalization, and computational efficiency. Such models ensure either invariance (i.e. output representation remains unchanged regardless of input transformations) or equivariance (i.e. output representation transforms accordingly to input transformations), preserving relative spatial relationships.

Recent advancements in geometric deep learning have led to the development of rotation-equivariant[52–54] and permutation-equivariant graph neural networks,[55,56] which explicitly encode spatial symmetries, improving PRL.[31,32,35,39] By inherently preserving geometric properties, these models enhance the accuracy of predictions in tasks such as protein structure modeling and molecular interactions. Moreover, by eliminating the need for extensive data augmentation—such as generating rotated versions of protein structures—they reduce computational overhead while improving model robustness and efficiency.

In parallel, invariant models address spatial symmetries by learning transformation-independent features, such as interatomic distances, bond angles, and backbone dihedral angles.[33,46,50,57] These geometric descriptors remain constant under rigid-body transformations, ensuring that identical protein structures in different orientations yield the same representation. For instance, atomic distances, bond angles, and dihedral angles capture essential spatial relationships while being inherently independent of absolute positioning. This approach simplifies learning, reduces computational complexity, and enhances generalization across diverse protein structures. By eliminating orientation and positional variance, invariant models streamline training while maintaining high

predictive performance.

# 5. Multimodal-Based Approaches

Multimodal-based approaches aim to integrate multiple sources of protein data, such as sequences, structures, and functional annotations, to create richer and more comprehensive representations. Each modality captures distinct aspects of a protein's characteristics: sequences provide information on the amino acid composition and evolutionary history, structures reveal the three-dimensional conformation critical for function, and functional annotations describe biological roles, molecular interactions, and phenotypic effects. By combining these complementary data types, multimodal approaches address the limitations inherent in single-modality models, leading to more accurate, robust, and generalizable protein representations.

## 5.1. Integration of Sequence and Structure Representations

One of the most prominent directions in multimodal PRL involves the integration of sequence and structure data. While Protein Language Models (PLMs) have demonstrated remarkable success in capturing evolutionary and contextual information from large-scale sequence datasets, their capability to encode the three-dimensional structural context essential for protein function remains limited. This limitation arises because sequence-only PLMs primarily model linear dependencies and lack explicit information about the spatial arrangement of residues. To address this gap, researchers have developed three primary strategies for incorporating structural information into PLMs to enhance their ability to model protein functions and interactions:

1. The first approach focuses on integrating sequence-derived information into structure-based models. In this method, PLMs are used to extract amino acid embeddings from sequences, which are then treated as node features in residue-level graphs.[58–62] These graphs are constructed using structural data, capturing spatial dependencies between residues through edges defined by proximity, interaction types, or geometric constraints. This strategy leverages

the contextual richness of PLMs while grounding the representations in the protein's three-dimensional structure, resulting in more comprehensive and informative models.

2. The second approach takes the reverse direction by integrating structural information directly into the training of PLMs. Instead of relying solely on sequence data, this method embeds geometric and spatial properties into the PLM framework, allowing the model to process structural features alongside amino acid sequences.[63–67] By incorporating structural cues, such as residue spatial proximity or geometric constraints, PLMs can generate richer, more biologically relevant representations that capture both the linear sequence and three-dimensional conformation of proteins.

3. The third approach treats sequence and structural data as distinct modalities, encoding them separately using models tailored to each.[42,57,68–73] The resulting representations are then fused through techniques such as attention mechanisms, contrastive learning, or embedding concatenation. This multi-branch strategy allows models to learn complementary features from both modalities while preserving the unique contributions of each, leading to enhanced representation quality and generalizability.

## 5.2. Integration of Sequence and Functional Representations

While Protein Language Models (PLMs) have advanced the understanding of sequence patterns, they often lack the explicit biological knowledge required to accurately capture protein function. Sequence data alone provides limited insight into a protein's specific roles or its molecular interactions within cellular environments, which can hinder model performance in tasks such as protein function prediction, protein contact mapping, and protein-protein interaction analysis. Given that a protein's structure fundamentally determines its function, models can be significantly improved by incorporating functional knowledge that highlights similarities among proteins with comparable shapes or sequences. Functional annotations provide comprehensive insights into a protein's biological processes, molecular functions, and cellular components. Integrating such functional data into PLMs enriches protein representations, allowing models to more effectively generalize

14

across diverse biological tasks and applications.

To address the limitations of sequence-only models, functional data is typically treated as a separate modality. Researchers employ biomedical language models like BioBERT[74] and PubMedBERT,[75] which are trained on large-scale biomedical literature and ontologies, to generate functional embeddings that capture rich contextual and domain-specific knowledge. These functional embeddings are then fused with sequence-derived representations using methods such as contrastive learning or joint embedding frameworks.[76–79] This multimodal integration enhances the representational capacity of PLMs, enabling more accurate and biologically informed predictions across a range of functional and interaction-based tasks.

## 5.3. Unified Representations of Sequence, Structure, and Function

Unified representations aim to simultaneously integrate sequence, structural, and functional data into a cohesive modeling framework, providing a comprehensive view of protein characteristics. While sequence-based models capture evolutionary and contextual relationships, and structure-based models provide insights into three-dimensional conformation, integrating functional annotations adds critical biological context regarding a protein's roles and interactions within cellular systems. By combining these modalities, unified representations offer a holistic understanding of proteins, enabling models to capture the complex interplay between a protein's linear sequence, its folded structure, and its biological function.

To achieve this, researchers typically adopt multi-branch architectures, where each modality—sequence, structure, and function—is encoded separately using specialized models. The outputs from these models are then fused using techniques such as attention mechanisms, contrastive learning, or cross-modal embedding fusion.[80,81] This approach ensures that the unique contributions of each modality are preserved while enabling the model to learn complementary features that enhance overall representation quality. By capturing the full complexity of protein biology, these multimodal approaches pave the way for more accurate, generalizable, and biologically meaningful models.

# 6. Complex-Based Approaches

Complex-based approaches focus on learning representations of entire protein complexes, capturing the structural and functional relationships between interacting molecules. These methods model interactions between proteins and ligands or protein-protein assemblies, which are fundamental to biological processes such as enzyme activity, signal transduction, and molecular recognition. By explicitly representing protein complexes, these approaches aim to learn context-aware embeddings that encode both individual molecular properties and interaction-specific features, facilitating tasks such as binding affinity prediction, molecular docking, and complex structure generation.

## 6.1. Protein-Ligand Complex Representation

Protein-ligand complex representation aims to learn a unified representation that effectively captures the interactions between proteins and small molecules, enabling models to better understand binding mechanisms and affinity relationships. Rather than representing the entire protein structure, these approaches typically focus on the binding pocket, the localized region where molecular interactions occur. By concentrating on this functionally relevant site, models can more accurately capture the spatial, chemical, and sequence-dependent features that govern ligand binding and molecular recognition. Approaches for protein-ligand complex representation can be broadly classified into two following categories, differentiated by their use of 3D structural information.

The first category follows a dual-encoder framework,[82–88] where proteins and ligands are encoded separately, without explicitly modeling their spatial interactions. Protein representations are commonly derived from protein language models (PLMs) or graph-based encoders, while ligand representations are generated using SMILES-based transformers, molecular fingerprints, or graph neural networks (GNNs). Since these models do not require pre-existing 3D structural data for the full complex, they are particularly advantageous for high-throughput virtual screening and cases where structural information is unavailable. During fine-tuning, a lightweight interaction module

16

is introduced to align and refine these representations for downstream tasks such as binding affinity prediction and molecular docking. This modular approach enhances scalability and facilitates transfer learning across diverse protein-ligand datasets.

The second category focuses on fine-grained interaction modeling, where 3D structural information of the full protein-ligand complex is explicitly incorporated to learn atomic or residue-level interactions. Unlike dual-encoder approaches, which align representations post hoc, these methods directly encode spatial and chemical dependencies within the binding pocket, capturing fine-grained molecular interactions. The binding mechanism of protein-ligand complexes is highly intricate, involving a diverse array of non-covalent forces such as $\pi$-stacking, $\pi$-cation interactions, salt bridges, water bridges, hydrogen bonds, hydrophobic interactions, and halogen bonds.[89] Previous studies[90–92] have demonstrated that explicitly modeling these atomic-level forces significantly improves predictive performance, underscoring the importance of 3D molecular force modeling in PRL.

## 6.2. Protein Complex Representation

Protein complexes, consisting of two or more interacting polypeptide chains, are fundamental to biological processes such as signal transduction, immune response, enzymatic activity, and molecular assembly. These complexes exhibit significant diversity in size, stability, and interaction mechanisms, ranging from small, transient interactions to large, stable macromolecular structures such as ribosomes, chaperones, and antibody-antigen complexes. The representation of protein complexes closely follows the methodology used in protein-ligand complex modeling, where the goal is to learn a joint representation that captures both individual protein properties and interaction-specific features. Approaches to this problem can be broadly categorized based on whether they explicitly encode the full complex structure or rely on separate representations of individual proteins.

Structure-free approaches, which encode individual proteins separately before aligning them in a shared representation space, are widely applied to various challenges in protein complex mod-

eling.[93–95] These approaches do not explicitly model the structural constraints of the full complex during encoding. Instead, interactions are inferred post hoc through mechanisms such as contrastive learning or learned scoring functions. Among the many tasks tackled by these approaches, protein complex structure prediction, determining the spatial arrangement of interacting polypeptides given only sequence or monomeric structural information, is one of the most critical.[96–98]

Structure-aware approaches, on the other hand, explicitly encode the entire 3D structure of the protein complex, incorporating spatial and physicochemical constraints to model binding interfaces and interaction dependencies.[99–101] These methods leverage geometric deep learning, graph-based modeling, and energy-based optimization techniques to refine complex formation by capturing atomic-level interactions and physical constraints. Because these approaches explicitly encode the joint structural context of interacting proteins, they can provide more interpretable predictions, improve docking accuracy, and enhance the design of protein-protein interactions for therapeutic applications. However, they require high-quality structural data, which may not always be available, and must balance computational efficiency with accuracy when handling large macromolecular assemblies.

# 7. Databases for Protein Representation Learning

Large-scale databases are the backbone of PRL, providing the essential sequence, structure, and functional data required for model training, validation, and benchmarking. This section introduces key databases widely used in PRL, categorized into sequence, structure, and protein function databases, as summarized in Table 2. These databases serve as fundamental resources for developing and evaluating protein representation models, enabling advancements in various biological and computational applications.

Table 1: Comparative Summary of Strengths and Limitations of PRL Approaches.

| PRL Approach | Strengths | Limitations |
|---|---|---|
| **Feature-Based** | • Simple and interpretable.<br>• Leverages well-established biochemical properties.<br>• Compatible with traditional machine learning models. | • Requires manual feature selection.<br>• May oversimplify protein characteristics.<br>• Struggles to generalize across diverse protein-related tasks. |
| **Sequence-Based** | • Utilizes large-scale protein sequence data for training Protein Language Models (PLMs).<br>• Captures evolutionary and contextual sequence information. | • MSA-based approaches require homologous sequences and can be computationally expensive.<br>• Non-aligned methods may lack evolutionary context. |
| **Structure-Based** | • Encodes rich spatial and functional information.<br>• Symmetry-aware models improve generalization. | • Requires high-quality structural data.<br>• Computationally intensive, especially at the atomic level.<br>• Surface-based representations may be memory-intensive. |
| **Multimodal-Based** | • Integrates multiple data sources for comprehensive protein representations.<br>• Improves model accuracy and generalization. | • Requires large, high-quality multimodal datasets.<br>• Increases computational complexity.<br>• Challenges in effectively fusing heterogeneous data modalities. |
| **Complex-Based** | • Captures interaction-specific features in protein-ligand and protein-protein complexes. | • Requires high-quality structural data.<br>• Structure-free methods may lack spatial constraints. |

Table 2: Summary of widely used databases in Protein Representation Learning (PRL), categorized by sequence, structure, and function.

| Name | Data | Description |
|---|---|---|
| UniProtKB [102,103] | Sequence, Function | UniProtKB comprises two main components: Swiss-Prot, which contains manually curated protein sequences, and TrEMBL, which provides automatically annotated sequences from high-throughput studies. |
| UniRef [104] | Sequence | UniRef is a clustering system based on the UniProt protein database, reducing sequence redundancy and improving computational efficiency by grouping sequences at different identity thresholds. |
| Pfam [105] | Sequence | Pfam is a database of protein families and domains, where each family is defined by multiple sequence alignments and profile Hidden Markov Models (HMMs) to facilitate functional annotation. |
| MGnify [106] | Sequence | MGnify is a large-scale metagenomic database containing over 2.4 billion non-redundant protein sequences predicted from environmental microbiomes. |
| PDB [107] | Structure | The Protein Data Bank (PDB) is a repository of experimentally determined biomolecular structures, obtained through techniques such as X-ray crystallography, NMR spectroscopy, and cryo-electron microscopy. |
| AlphaFold DB [108] | Structure | AlphaFold DB stores protein structures predicted using the AlphaFold series, offering high-confidence computational models that expand structural coverage beyond experimentally resolved proteins. |
| ESMAtlas [108] | Structure | ESMAtlas is a large-scale protein structure prediction database generated using ESMFold, containing over 617 million predicted structures, including millions of novel proteins not yet experimentally characterized. |
| Gene Ontology [109] | Function | The Gene Ontology (GO) database provides structured annotations of protein functions, categorizing them into biological processes, molecular functions, and cellular components to support functional characterization. |

# 8. Applications

Protein Representation Learning (PRL) has revolutionized biological research and biomedical innovation, providing data-driven solutions for understanding and engineering proteins. By learning meaningful representations of protein sequences, structures, and interactions, PRL enables a wide range of applications, from predicting fundamental protein properties to advancing drug discovery. As illustrated in Fig. 4, the applications of PRL can be broadly categorized into four main domains: protein property prediction, protein structure prediction, protein design and optimization, and drug discovery.

In the following sections, each application area is presented through a structured discussion that includes: (i) a problem statement defining the task, (ii) an explanation of its significance in biological and biomedical contexts, (iii) an overview of commonly used benchmarks for evaluating model performance, and (iv) a discussion of models that have advanced the field, with observations on their representation learning strategies. This structure provides a comprehensive and comparative perspective on how PRL is shaping advancements in protein science and biopharmaceutical applications.

## 8.1. Protein Property Prediction

Predicting protein properties involves estimating key characteristics such as solubility, stability, enzymatic activity, and binding affinity from sequence or structural data. These properties are fundamental to protein function and have significant implications in biotechnology, medicine, and drug discovery. Accurate prediction models reduce experimental costs and accelerate the design of novel therapeutics and industrial enzymes, making protein property prediction a central task in PRL.

Benchmarks for evaluating protein property prediction models are categorized based on the type of property being assessed. For individual protein properties, Tasks Assessing Protein Embeddings (TAPE)[110] provides a standardized evaluation framework, including fluorescence predic-
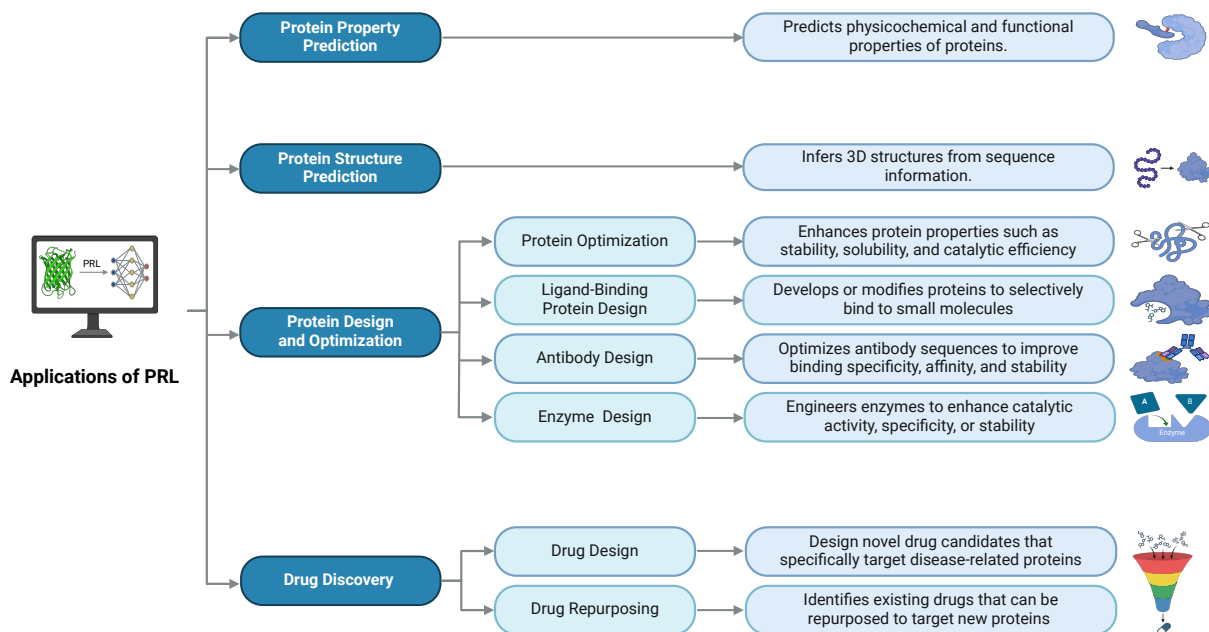
Figure 4: Overview of key applications of Protein Representation Learning (PRL). PRL enables advancements in multiple domains, including protein property prediction, structure prediction, protein design and optimization, and drug discovery.

tion and protease stability assessment. Another widely used benchmark, UniProt,[102] serves as a comprehensive database of protein sequences, functional annotations, and metadata from diverse organisms, offering a valuable resource for assessing sequence-based property predictions across various protein families and functional contexts. In contrast, interaction-based properties—such as protein-ligand and protein-protein interactions—require different benchmarks that incorporate structural and binding information. Prominent datasets in this category include PDBBind,[111] Leak Proof PDBBind,[112] Davis,[113] KIBA,[114] PLINDER[115] and PINDER,[116] which evaluate models based on binding affinity predictions. These datasets serve as widely recognized standards, enabling consistent model comparison and performance assessment across different protein property prediction tasks.

With the remarkable advancements in Protein Language Models (PLMs), most current research leverages these models to extract meaningful representations of protein sequences, applying them to a wide range of property prediction tasks. These models have been successfully employed for

both individual protein properties [7,117,118] as well as interaction-based properties. [93,119,120] Beyond sequence-based representations, many efforts have explored integrating multiple modalities, such as structural, physicochemical, and functional data, to further enhance protein representations beyond what PLMs alone can achieve. [42,58,72,121] By incorporating additional biological information, these approaches aim to provide a more comprehensive view of protein function and interactions, particularly in tasks where spatial and contextual features play a significant role, such as protein-ligand binding and protein-protein interactions.

## 8.2. Protein Structure Prediction

Protein structure prediction, also known as protein folding, aims to determine a protein's three-dimensional conformation from its amino acid sequence. Since protein structure dictates function, accurately predicting spatial arrangements is essential for understanding molecular interactions, drug binding, and enzymatic activity. Traditionally, protein structures have been determined using experimental techniques such as X-ray crystallography [122] and cryo-electron microscopy (Cryo-EM), [123] which provide high-resolution structural insights. However, these methods are often time-consuming, costly, and limited by challenges such as crystallization feasibility and sample preparation. [124] Advances in computational structure prediction have transformed the field by providing faster, scalable alternatives that complement experimental approaches, broadening the scope of protein engineering and biomedical research.

To evaluate the accuracy of structure prediction models, several benchmarking platforms have been established. The Critical Assessment of Techniques for Protein Structure Prediction (CASP), [125] a biennial blind assessment, is widely regarded as the gold standard for evaluating protein folding models. In addition, the Continuous Automated Model EvaluatiOn (CAMEO) [126] provides real-time assessments by continuously comparing predictions against newly released experimental structures. These benchmarks are instrumental in validating and refining computational models, ensuring their reliability and applicability in structural biology.

Protein structure prediction models can be broadly categorized into MSA-based and MSA-free

approaches. MSA-based methods leverage multiple sequence alignments to capture evolutionary relationships, often resulting in higher accuracy for proteins with rich evolutionary data. Notable models in this category include AlphaFold,[8] RosettaFold,[127] OpenFold,[128] and MSAGPT.[11] In contrast, MSA-free methods predict structures directly from single sequences without relying on evolutionary context, making them more efficient and applicable to proteins with limited sequence homology. Examples of MSA-free models include ESMFold,[13] OmegaFold,[129] and HelixFold-Single.[130] While these models have greatly improved protein structure prediction, they are primarily limited to static conformations, as they are trained on crystallographic and experimentally resolved structures. However, proteins are inherently dynamic, undergoing conformational changes essential to their function. To address this limitation, recent advancements have introduced generative models capable of predicting conformational ensembles, such as BioEmu.[131]

## 8.3. Protein Design and Optimization

Protein design and optimization involve strategies for engineering proteins with specific structural and functional properties. Broadly, these strategies fall into two main categories: sequence-based design, which optimizes amino acid sequences to achieve target properties, and structure-based design, which leverages protein structural features—either by generating backbones and applying inverse folding to identify compatible sequences or by co-designing sequence and structure to enhance functional and structural outcomes. These approaches have been applied to various protein types, including individual proteins with optimized functionality, ligand-binding proteins for molecular interactions, antibodies for therapeutic applications, and enzymes for improved catalytic efficiency in industrial and biomedical settings.

### 8.3.1. Individual Protein Property Optimization

Individual protein property optimization focuses on enhancing a protein's intrinsic characteristics, such as stability, activity, specificity, and catalytic efficiency, without explicitly considering interactions with other molecules. Traditional approaches, such as directed evolution with random

mutagenesis, [132–134] iteratively introduce mutations and select for improved variants. While these methods have been widely used and have led to many successful optimizations, they often require extensive experimental screening. Computational strategies have been developed to facilitate a more systematic exploration of sequence space, offering alternatives that reduce experimental reliance while aiming to improve efficiency in protein engineering.

To evaluate optimization strategies, several benchmark datasets have been established for sequence-based optimization, each representing distinct challenges across protein targets, originating organisms, sequence lengths, dataset sizes, and fitness objectives. Notable benchmarks include Green Fluorescent Protein (avGFP), [135] Adeno-Associated Viruses (AAV), [136] TEM-1 $\beta$-Lactamase (TEM), [137] Ubiquitination Factor Ube4b (E4B), [138] Aliphatic Amide Hydrolase (AMIE), [139] Levoglucosan Kinase (LGK), [140] Poly(A)-binding Protein (Pab1), [141] and SUMO E2 Conjugase (UBE2I). [142] These datasets enable standardized assessment of different protein optimization methods.

Earlier approaches focused on optimizing protein in discrete sequence space. These methods explored mutation-based search strategies, reinforcement learning, and PLMs to propose beneficial sequence modifications. [143,144] Building upon this foundation, recent advancements in PRL have introduced computational frameworks capable of mapping protein sequences into fitness landscapes. These models aim to capture sequence-function relationships and facilitate direct sequence optimization based on learned representations. Unlike experimental screening, PRL-based approaches rely on large-scale protein datasets to infer meaningful sequence features and predict fitness variations. A common strategy in fitness landscape modeling integrates sequence and fitness information within a machine-learning framework. One such approach involves jointly training an autoencoder with a prediction network, where an encoder extracts sequence features, a decoder reconstructs sequences, and a separate model predicts fitness values from the latent space. [145,146] Protein Language Models (PLMs) are often used as encoders to capture contextual sequence representations, and the prediction network guides latent space organization based on fitness variations.

Other studies have explored methods to improve fitness landscape modeling through smoothing techniques, aiming to enhance search efficiency in sequence space. Graph-based smoothing

techniques have been introduced to incorporate protein similarity graphs, promoting structural continuity in the fitness landscape.[147,148] These methods aim to reduce abrupt fitness changes caused by minor sequence variations, improving the robustness of optimization algorithms. Overall, while experimental and computational strategies each offer advantages and limitations, ongoing research continues to explore ways to enhance protein optimization efficiency while ensuring the reliability and interpretability of computational predictions.

### 8.3.2. Ligand-Binding Protein Design

Ligand-binding protein design involves engineering proteins to interact with small molecules such as drugs, metabolites, or signaling compounds. These interactions are central to biological processes, making ligand-binding protein design essential for applications in biosensors, precision medicine, and enzyme-mediated synthesis. Achieving high specificity and affinity for target ligands requires a detailed understanding of binding site geometry, molecular interactions, and binding thermodynamics, which computational methods have been increasingly employed to model and optimize.

Several benchmark datasets support the evaluation of ligand-binding protein design strategies by providing experimentally validated protein-ligand complexes. PDBbind[111] offers a curated collection of protein-ligand structures with binding affinity measurements, serving as a reference for affinity prediction and binding site modeling. CSAR (Community Structure Activity Resource)[149] provides high-resolution crystallographic data with a diverse range of ligand complexes, offering a structured framework for evaluating binding specificity across various protein targets. These datasets are widely used to assess structure-based and sequence-based ligand-binding design approaches.

Computational strategies for ligand-binding protein design typically leverage structural insights from native protein-ligand complexes to refine side-chain interactions and backbone conformations, with the aim of enhancing binding affinity.[150–153] Many approaches follow a structure-first paradigm, in which protein backbones are generated first, followed by inverse folding tech-

26

niques[154,155] to identify sequences that can adopt predefined structural conformations. This prioritization of structure prediction before sequence determination allows for the design of proteins that conform to specific ligand-binding geometries and can also aid in addressing cases where binding sites are not well defined.

Despite progress, challenges remain, particularly in the structural validation of designed binding modes. While computational docking has been applied to generate novel binders by modifying scaffolds and loop geometries, the accuracy and stability of predicted interactions often require high-resolution experimental confirmation. Additionally, many design strategies emphasize a limited set of hotspot residues for scaffold placement, which may restrict the exploration of diverse interaction modes, particularly for targets lacking well-defined pockets or binding clefts.[156,157] To address these limitations, alternative approaches have investigated ligand-binding protein design based on sequence information alone, bypassing the reliance on structural data.[87] These models leverage Protein Language Models (PLMs) to extract functional sequence motifs associated with ligand binding, allowing for the design of proteins with desired binding properties even when structural information is unavailable or difficult to obtain.

### 8.3.3. Antibody Design

Antibodies are specialized immune proteins that recognize and bind to specific antigens, such as pathogens, toxins, or diseased cells, marking them for neutralization or destruction. As central components of adaptive immunity, antibodies exhibit high specificity and strong binding affinity, enabling precise immune responses. Their function is primarily dictated by variable regions, particularly the complementarity-determining regions (CDRs), which mediate antigen recognition. Antibody design focuses on engineering antibodies with improved affinity, specificity, and stability, with broad therapeutic applications in cancer, autoimmune disorders, and infectious diseases.

Benchmark datasets play a critical role in evaluating sequence-based and structure-based antibody design. OAS (Observed Antibody Space)[158] provides over 2 billion immune repertoire sequences from various immune states, species, and individuals, supporting sequence-based anti-

27

body modeling. SAbDab (Structural Antibody Database)[159] compiles annotated antibody structures from the PDB, including affinity data and sequence annotations, serving as a key resource for structure-based antibody modeling. AbBind[160] includes 1,101 antibody-antigen mutants across 32 complexes with experimentally measured binding free energy changes, allowing for the study of affinity variations upon mutation. These datasets support the evaluation and development of antibody design approaches, ensuring standardized model assessment across different strategies.

Computational strategies for structure-based antibody design primarily focus on CDRH3 loops, which are highly variable and critical for antigen binding. Two main methodologies have emerged: (i) 3D backbone generation, which designs structurally realistic CDRH3 loops for stable antigen recognition,[161] and (ii) binding affinity prediction, which models how sequence modifications impact structural stability and antigen binding energy.[162] While structure-based approaches have been widely used, a major challenge is the limited availability of high-resolution antibody structures, which constrains direct structural modeling. To address this, sequence-based models leverage Protein Language Models (PLMs) trained on large-scale immune repertoires to learn patterns associated with high-affinity antibodies. These methods enable de novo antibody generation and affinity maturation without requiring structural information.[163] Additionally, sequence-based antibody design can be framed as an optimization problem, where binding affinity serves as the fitness metric, guiding the search for improved variants—similar to optimization strategies used for individual protein properties.[164,165] A promising direction in antibody design is the integration of sequence and structural information, leading to sequence-structure co-design approaches.[166–170] By combining both modalities, these models aim to capture relationships between sequence variability and structural adaptation, facilitating the precise engineering of antibody specificity, stability, and affinity. This hybrid framework expands the design space beyond naturally occurring immune repertoires, potentially accelerating antibody discovery and improving antigen recognition strategies.

### 8.3.4. Enzyme Design

Enzymes are biological catalysts that accelerate chemical reactions by lowering activation energy, playing fundamental roles in metabolism, signal transduction, and biomolecule synthesis and degradation. Their high specificity and efficiency make them indispensable in both cellular processes and industrial applications, including pharmaceutical production, biofuel synthesis, food processing, and environmental remediation. Enzyme design focuses on engineering or optimizing enzymes to enhance catalytic efficiency, stability, and substrate specificity for targeted reactions. Traditional approaches such as directed evolution and rational design have been widely applied but often require extensive experimental screening.[171] Computational strategies now enable in silico enzyme prediction and optimization, facilitating a more efficient design process.

Benchmark datasets play a critical role in assessing enzyme design models by providing standardized evaluation frameworks for sequence-function relationships and catalytic performance. UniProt[102] offers a comprehensive database of protein sequences and functional annotations across diverse organisms, supporting sequence-based enzyme design and optimization. BRENDA[172] compiles detailed enzyme data, including sequence, structure, and kinetic parameters, enabling validation of catalytic efficiency, stability, and substrate specificity in designed enzymes. These benchmarks provide essential resources for the evaluation and development of computational enzyme engineering strategies.

Recent advancements in enzyme design for specific Enzyme Commission (EC) classes have demonstrated promising results, with models generating sequences that closely resemble reference enzymes while maintaining desired catalytic functions.[28,173] Beyond EC-based classification, alternative approaches have explored de novo enzyme generation by assembling active site and scaffold libraries, followed by refinement algorithms to improve functionality.[174] While EC classification provides a structured framework, relying solely on predefined categories may limit model generalization to novel, unseen reactions. To address these challenges, recent research has shifted toward reaction-conditioned enzyme design, which directly models enzyme-substrate relationships rather than relying on predefined EC classifications.[175,176] This approach enables greater

flexibility in enzyme generation, allowing models to learn catalytic patterns from reaction-specific data.

Another key challenge in enzyme sequence design is the limited understanding of enzyme-substrate catalytic mechanisms. Even when designed enzyme sequences fold correctly into 3D structures, catalytic pocket formation and binding interactions often remain poorly characterized. Recent efforts have integrated generative models for enzyme scaffolds, active sites, and protein language models to improve enzyme-substrate interaction modeling.[177] By generating enzyme-substrate binding structures, these methods aim to refine catalytic mechanism prediction and support the design of enzymes capable of catalyzing novel reactions.

## 8.4. Drug Discovery

Drug discovery utilizes protein structures and sequences to develop or repurpose small molecules that bind with high specificity and affinity. By analyzing the spatial, physicochemical, and sequence properties of protein targets, researchers aim to design drugs that minimize off-target effects and enhance therapeutic efficacy. Drug discovery is particularly critical for diseases with well-characterized protein targets, including cancer, infectious diseases, and neurodegenerative disorders. This section covers two key tasks: drug design, which focuses on generating novel therapeutic compounds, and drug repurposing, which identifies new applications for existing drugs.

### 8.4.1. Drug Design

Drug design focuses on generating novel small molecules tailored to bind specific protein targets with optimized affinity, stability, and pharmacokinetics. By leveraging both structural and sequence-based information, drug design aims to maximize specificity while minimizing off-target effects, making it essential for treating diseases where precise molecular interactions are crucial, such as cancer and neurodegenerative disorders.

To support model development and evaluation, benchmark datasets provide standardized training and testing resources. CrossDocked[178] consists of cross-docked protein-ligand pairs, chal-

lenging models to generalize across diverse binding poses, a crucial ability for designing molecules across multiple protein targets. Binding MOAD (Mother of All Databases)[179] offers a large collection of high-resolution protein-ligand complexes with experimentally measured binding affinities, serving as a foundation for refining docking precision and affinity predictions. Together, these datasets facilitate robust evaluation of structure-based drug design (SBDD) models.

Recent advances in PRL have led to deep learning models that generate molecular structures conditioned on 3D representations of protein binding pockets. These models capture the geometric and physicochemical attributes of binding sites, facilitating the design of small molecules with high affinity and specificity. A predominant strategy in structure-based drug design (SBDD), this pocket-conditioned molecular generation paradigm includes notable models such as Pocket2Mol,[180] ResGen,[181] PocketFlow,[182] DeepICL,[183] and DiffSBDD,[184] all of which utilize 3D structural data to guide molecular design. By anchoring molecular generation to the spatial and chemical features of protein pockets, these approaches improve ligand design precision and streamline SBDD workflows.

Despite these advances, accurately defining binding pockets for novel protein targets remains a challenge, particularly when structural annotations are unavailable or ambiguous. While many models rely on predefined pockets, emerging efforts are shifting toward whole-protein drug design, allowing models to identify potential binding regions in a data-driven manner without requiring prior annotations.[72] This approach broadens the applicability of SBDD to less-characterized targets. Additionally, a complementary research direction focuses on sequence-based drug design, where molecular structures are generated directly from protein sequences, bypassing the need for 3D structural data.[185,186] These methods leverage Protein Language Models (PLMs) to extract sequence-level functional insights, offering an alternative strategy for designing drugs against targets with limited or unreliable structural information.

### 8.4.2. Drug Repurposing

Drug repurposing identifies new therapeutic applications for existing drugs, offering a faster and more cost-effective alternative to de novo drug development. This approach is particularly valuable for urgent medical needs, such as pandemics, where traditional drug discovery timelines may be impractical. Broadly, drug repurposing can be categorized into two primary strategies: (i) Virtual screening, which identifies existing drugs with high binding affinity to specific protein targets; and (ii) Drug-Target Interaction (DTI) prediction via link prediction, which models relational networks to uncover novel drug-protein associations.

To facilitate standardized evaluation, benchmark datasets provide essential validation resources for both virtual screening and DTI-based approaches. For virtual screening, DUD-E (Directory of Useful Decoys: Enhanced)[187] and LIT-PCBA[188] are widely used. DUD-E pairs protein targets with both active compounds and decoys, assessing a model's ability to differentiate true binders from non-binders in complex screening tasks. LIT-PCBA, derived from PubChem bioassays, offers experimentally validated active and inactive compounds, reflecting the challenges encountered in real-world drug repurposing efforts. For DTI-based link prediction, datasets such as PrimeKG,[189] Therapeutics Data Commons (TDC),[190] BindingDB,[191] and BioSNAP[192] provide extensive collections of validated drug-target interactions. These resources support the development and benchmarking of computational models that predict novel drug-target relationships, enabling systematic drug repurposing.

Virtual screening can be further categorized into docking-based and similarity-based methods. Docking-based approaches use molecular docking simulations, such as AutoDock Vina,[193] to computationally dock large drug libraries against protein targets and predict binding affinities. While effective, traditional docking is computationally intensive and limits scalability for high-throughput screening. To overcome these limitations, deep learning-based docking models, including P2Rank,[194] EquiBind,[195] TANKBind,[196] DiffDock[197] and HelixDock,[198] have been developed to predict binding poses and affinities with reduced computational costs. These models leverage 3D protein-ligand complexes to learn structural binding interactions and improve dock-

ing accuracy. In contrast, similarity-based approaches rely on the "similarity principle", which assumes that structurally similar molecules exhibit similar biological properties.[199] Rather than explicitly docking compounds, these methods use learned molecular representations to rapidly identify structurally related drug candidates. Recent advances in PRL have enabled drug and protein embeddings to be mapped into a shared representation space, facilitating efficient similarity searches. Unlike docking-based methods that depend on 3D structural data, similarity-based approaches do not require explicit protein-ligand modeling. Instead, they employ pretrained models, such as Protein Language Models (PLMs) and Molecular Graph Neural Networks (GNNs), to encode proteins and ligands separately. Notable examples include ConPLex,[200] DrugCLIP,[201] and SPRINT,[202] which use contrastive learning to map proteins and drugs into a joint representation space, allowing for rapid and scalable virtual screening without the computational cost of docking simulations.

Beyond virtual screening, DTI-based link prediction offers an alternative paradigm for drug repurposing by formulating drug-target interactions as a graph-based link prediction task. In this approach, drugs and proteins are represented as nodes, while known interactions form edges. A key challenge in DTI-based models is how drug and protein nodes are represented, as the quality of node embeddings directly impacts model performance. Recent methods have introduced various node embedding strategies, which can be classified into sequence-based and multimodal embeddings. Sequence-based embeddings use PLMs for protein sequences and SMILES-based transformers for molecular representations,[203,204] encoding drugs and proteins based on their primary sequences without requiring structural data. Meanwhile, multimodal embeddings integrate diverse biological information sources, combining sequence, structure, and auxiliary biological data, such as functional annotations, side-effect similarities, and pathway interactions, to construct more comprehensive representations.[205,206] By capturing both structural and functional relationships, these models improve the predictive accuracy of link prediction tasks, making DTI-based drug repurposing a valuable complement to virtual screening approaches by providing insights into drug-target associations that may not be captured through traditional binding-based models.

33

# 9. Future Directions and Open Challenges

While PRL has made significant strides, several challenges and opportunities remain that could shape the future development of the field. This section explores four key directions that we believe will drive further advancements: extending PRL methodologies to other biological sequences such as DNA and RNA to uncover regulatory mechanisms and gene functions, improving model explainability to enhance interpretability and trust, scaling PRL models to improve efficiency and accessibility, and enhancing generalization to improve robustness across unseen and mutated proteins, as summarized in Fig. 5.

## 9.1. Biological Sequence Representation Learning

After significant advancements in Protein Representation Learning (PRL), a natural progression is the extension of these methods to other biological macromolecules, such as DNA and RNA. Learning meaningful representations of DNA and RNA sequences is critical for understanding gene regulation, epigenetics, alternative splicing, and disease mechanisms. By applying methodologies developed in PRL to these biological sequences, researchers can leverage the power of deep learning models to extract insights that may otherwise remain hidden within complex genomic and transcriptomic data.

The representation learning of DNA and RNA shares foundational similarities with protein representation learning, as both involve extracting meaningful features from biological sequences to predict function, interactions, and structural properties. Like proteins, DNA and RNA sequences follow a linear arrangement of biological "letters" that encode essential biological information. Advances in self-supervised learning and transformer-based architectures, which have successfully modeled protein sequences, can similarly be applied to DNA and RNA to capture sequence patterns, evolutionary conservation, and functional motifs. Recent efforts to develop genome language models have shown promising results, demonstrating their potential in learning meaningful representations of genomic sequences and improving the prediction of gene function, regulatory

34

elements, and chromatin accessibility.[207–209]

However, DNA and RNA present unique challenges distinct from proteins. Unlike proteins, which have a well-defined alphabet of 20 amino acids and a structured hierarchical organization, DNA and RNA sequences consist of only four nucleotides (A, T/U, G, and C) but often span much longer sequences, reaching tens of thousands or even millions of base pairs. This length disparity makes learning effective representations computationally intensive and necessitates efficient sequence compression strategies. Additionally, while proteins fold into stable three-dimensional structures that largely dictate function, DNA and RNA often exhibit dynamic secondary and tertiary structures influenced by sequence context, environmental factors, and cellular conditions. Capturing these structural variations requires novel modeling approaches that integrate both sequence-based and structural information. Furthermore, the function of DNA and RNA is heavily modulated by epigenetic modifications, chromatin accessibility, and interactions with regulatory proteins. Unlike proteins, where function is often inferred from sequence and structure, the functional state of a genomic region is highly context-dependent, requiring models to account for regulatory mechanisms beyond the primary sequence.

## 9.2. Scaling PRL Models for Large-Scale Applications

As biological datasets continue to grow in size and complexity, the scalability of PRL models remains a significant challenge. Many existing methods struggle with large-scale protein representations, requiring substantial computational resources that limit accessibility and real-world applicability. Addressing these challenges requires advancements in model efficiency and optimization techniques.

One promising direction for improving scalability is model distillation, a technique that has seen significant progress in Natural Language Processing (NLP) and Computer Vision (CV). In NLP, approaches such as DistilBERT[210] and TinyBERT[211] have successfully compressed large transformer models while preserving much of their predictive power. Similarly, in CV, TinyViT[212] and DearKD[213] leverage distillation to create more compact and computationally efficient models.

35

These techniques demonstrate the potential of distilling knowledge from large, computationally expensive models into smaller, more efficient variants without significant loss of accuracy. Recent efforts have explored applying model distillation techniques to Protein Representation Learning (PRL) to improve scalability and efficiency. For instance, OpenFold[128] employs knowledge distillation to create a more computationally efficient version of AlphaFold.

Applying model distillation to PRL could enhance efficiency by enabling smaller models to retain essential biological insights while significantly reducing memory and computational overhead. This would be particularly beneficial for large-scale virtual screening, protein design, and low-data scenarios where fine-tuning massive models is impractical. Additionally, integrating distillation with techniques such as low-rank adaptations, mixed-precision training, and distributed learning could further enhance efficiency while maintaining high performance.

## 9.3. Improving Generalization of PRL Models

While scalability focuses on computational efficiency, generalization remains a critical challenge in PRL. Many models struggle to generalize to unseen proteins, such as those from novel species, synthetic biology applications, or de novo protein designs, limiting their applicability in biomedical and biotechnological research. Beyond handling completely novel proteins, PRL models must also account for mutated proteins, where small amino acid substitutions, deletions, or insertions can significantly alter a protein's function, stability, or interactions. Understanding how mutations affect learned representations is crucial for applications such as disease mutation analysis, where genetic variants influence protein function; protein engineering, where specific mutations can optimize protein stability and activity; and drug resistance prediction, where mutations in pathogens impact therapeutic efficacy. However, the availability of labeled data for both novel proteins and mutant proteins is often extremely limited, making it difficult to train models that can reliably predict their properties.

To address this, zero-shot and few-shot learning approaches need to be developed, enabling PRL models to generalize effectively to unseen proteins and mutations with minimal supervi-

sion. Recent studies have demonstrated promising results in this area, showing that pretrained PRL models can infer functional properties and structural effects of mutations even in data-scarce settings, improving predictive performance in mutation effect prediction and de novo protein engineering.[57,214]

## 9.4. Enhancing Explainability in PRL Models

Despite their success, many state-of-the-art PRL models operate as black boxes, making it difficult to interpret how they learn protein representations and make predictions. This lack of transparency hinders scientific discovery, limits trust in model outputs, and reduces their applicability in critical fields such as medicine, biotechnology, and drug development. As PRL models become increasingly complex—leveraging deep neural networks, large-scale pretraining, and self-supervised learning—there is a growing need for approaches that can make their predictions more interpretable and biologically meaningful.

One approach to enhancing explainability in PRL is the development of self-interpretable architectures that inherently capture biologically meaningful representations. Explainable AI techniques, such as attention mechanisms, saliency mapping, and feature attribution, allow models to highlight critical sequence motifs, structural domains, and evolutionary patterns that drive their predictions. By identifying key residues, structural elements, and physicochemical properties relevant to specific functions or interactions, these methods improve model transparency and reliability. This enhanced interpretability facilitates applications in functional protein annotation, biomarker discovery, and rational protein design. Recent studies have explored these approaches, applying various Explainable AI techniques to improve the interpretability of PRL models.[215,216]

Another promising direction is leveraging advances in Natural Language Processing (NLP) and Large Language Models (LLM) to make PRL models more accessible and interpretable. Inspired by recent successes in LLM, PRL models could be designed to generate natural language descriptions of protein functions, structural properties, and interactions, making predictions more intuitive for researchers in biology, medicine, and drug development. Beyond simple text-based outputs,

reasoning-based approaches could further enhance explainability by allowing models to logically infer protein functions and molecular interactions rather than relying solely on statistical pattern recognition. Techniques such as retrieval-augmented generation (RAG) and biomedical knowledge graphs (BKGs) could serve as foundational tools for reasoning, enabling models to retrieve relevant biological context from structured databases and scientific literature. Instead of merely generating descriptions, PRL models could use retrieved biological relationships and prior knowledge to construct logical justifications for their predictions, improving their transparency and reliability. Recent studies have already begun exploring these approaches, demonstrating promising results in enhancing model interpretability and providing biologically meaningful insights.[217,218] By integrating LLM-based reasoning, knowledge retrieval, and structured biological insights, PRL models could evolve into more explainable, evidence-driven systems, facilitating their adoption in research and clinical applications.

# 10. Conclusion

Protein representation learning (PRL) has emerged as a transformative approach to advancing our understanding of proteins and addressing critical challenges in molecular biology, biotechnology, and medicine. This review provides a comprehensive overview of PRL methodologies, categorized into five major approaches: feature-based, sequence-based, structure-based, multimodal, and complex-based. By analyzing the strengths and limitations of each approach, we offer insights to guide the selection of appropriate PRL strategies based on specific biological and computational requirements, ensuring their effective application across diverse research domains.

Beyond methodological advancements, we highlight key databases that underpin PRL research and discuss its diverse applications, including protein property prediction, structure modeling, protein design, and drug discovery. Despite the substantial progress in PRL, several open challenges remain, presenting opportunities for future research. Key directions for advancement include extending PRL methodologies to genomic representation learning, enhancing model scalability and
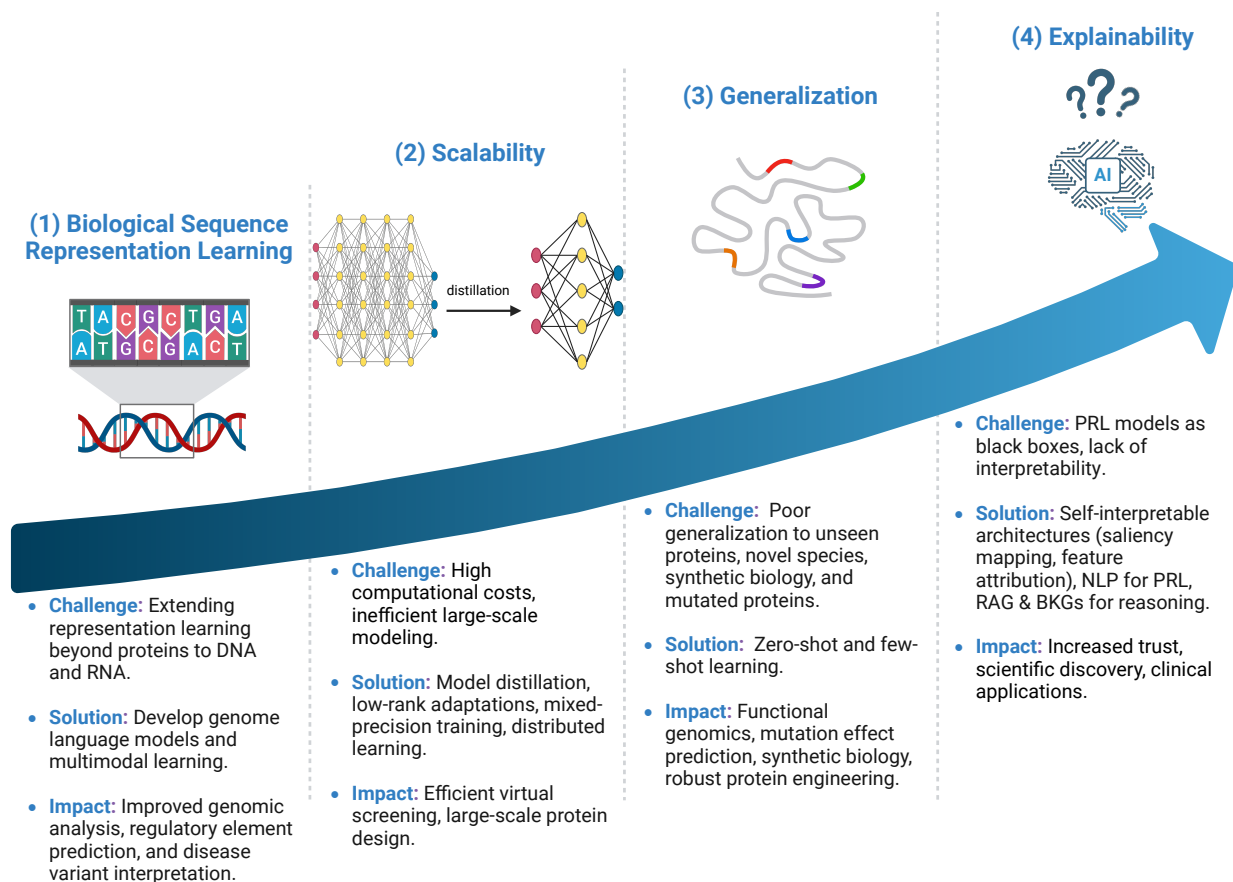
Figure 5: Key future directions in Protein Representation Learning (PRL). The figure outlines four critical challenges and potential advancements: (i) Expanding PRL to DNA and RNA representation learning to leverage shared methodologies while addressing unique challenges, (ii) Enhancing scalability to improve computational efficiency and accessibility of large-scale models, (iii) Strengthening generalization to ensure robustness across unseen proteins and genetic variations, and (iv) Advancing explainability to improve model interpretability and facilitate trust in biological and biomedical applications.

generalization, and improving interpretability to foster broader accessibility and adoption. Addressing these challenges will be crucial in refining PRL methodologies and expanding their impact across biological and biomedical research.

# References

(1)  Wu, L.; Huang, Y.; Lin, H.; Li, S. Z. A survey on protein representation learning: Retrospect and prospect. *arXiv preprint arXiv:2301.00813* **2022**,

(2) Xiao, Y.; Zhao, W.; Zhang, J.; Jin, Y.; Zhang, H.; Ren, Z.; Sun, R.; Wang, H.; Wan, G.; Lu, P., et al. Protein Large Language Models: A Comprehensive Survey. *arXiv preprint arXiv:2502.17504* **2025**,

(3) Heinzinger, M.; Rost, B. Teaching AI to speak protein. *Current Opinion in Structural Biology* **2025**, *91*, 102986.

(4) Bedbrook, C. N.; Yang, K. K.; Rice, A. J.; Gradinaru, V.; Arnold, F. H. Machine learning to design integral membrane channelrhodopsins for efficient eukaryotic expression and plasma membrane localization. *PLoS computational biology* **2017**, *13*, e1005786.

(5) Saladi, S. M.; Javed, N.; Müller, A.; Clemons, W. M. A statistical model for improved membrane protein expression using sequence-derived features. *Journal of Biological Chemistry* **2018**, *293*, 4913–4927.

(6) Kawashima, S.; Kanehisa, M. AAindex: amino acid index database. *Nucleic acids research* **2000**, *28*, 374–374.

(7) Yang, K. K.; Wu, Z.; Bedbrook, C. N.; Arnold, F. H. Learned protein embeddings for machine learning. *Bioinformatics* **2018**, *34*, 2642–2648.

(8) Jumper, J.; Evans, R.; Pritzel, A.; Green, T.; Figurnov, M.; Ronneberger, O.; Tunyasuvunakool, K.; Bates, R.; Žídek, A.; Potapenko, A., et al. Highly accurate protein structure prediction with AlphaFold. *nature* **2021**, *596*, 583–589.

(9) Rao, R. M.; Liu, J.; Verkuil, R.; Meier, J.; Canny, J.; Abbeel, P.; Sercu, T.; Rives, A. MSA transformer. International Conference on Machine Learning. 2021; pp 8844–8856.

(10) Zheng, W.; Wuyun, Q.; Li, Y.; Zhang, C.; Freddolino, P. L.; Zhang, Y. Improving deep learning protein monomer and complex structure prediction using DeepMSA2 with huge metagenomics data. *Nature Methods* **2024**, *21*, 279–289.

(11) Chen, B.; Bei, Z.; Cheng, X.; Li, P.; Tang, J.; Song, L. MSAGPT: Neural Prompting Protein Structure Prediction via MSA Generative Pre-Training. *arXiv preprint arXiv:2406.05347* **2024**,

(12) Meng, Q.; Guo, F.; Tang, J. Improved structure-related prediction for insufficient homologous proteins using MSA enhancement and pre-trained language model. *Briefings in Bioinformatics* **2023**, *24*, bbad217.

(13) Lin, Z.; Akin, H.; Rao, R.; Hie, B.; Zhu, Z.; Lu, W.; Smetanin, N.; Verkuil, R.; Kabeli, O.; Shmueli, Y., et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science* **2023**, *379*, 1123–1130.

(14) Hou, Q.; De Geest, P. F.; Vranken, W. F.; Heringa, J.; Feenstra, K. A. Seeing the trees through the forest: sequence-based homo-and heteromeric protein-protein interaction sites prediction using random forest. *Bioinformatics* **2017**, *33*, 1479–1487.

(15) Yao, Y.-h.; Lv, Y.-p.; Li, L.; Xu, H.-m.; Ji, B.-b.; Chen, J.; Li, C.; Liao, B.; Nan, X.-y. Protein sequence information extraction and subcellular localization prediction with gapped k-Mer method. *BMC bioinformatics* **2019**, *20*, 1–8.

(16) Sinai, S.; Kelsic, E.; Church, G.; Nowak, M. Variational auto-encoding of protein sequences. arXiv. 2017.

(17) Ding, X.; Zou, Z.; Brooks III, C. L. Deciphering protein evolution and fitness landscapes with latent space models. *Nature communications* **2019**, *10*, 5644.

(18) Almagro Armenteros, J. J.; Johansen, A. R.; Winther, O.; Nielsen, H. Language modelling for biological sequences–curated datasets and baselines. *BioRxiv* **2020**, 2020–03.

(19) Rives, A.; Meier, J.; Sercu, T.; Goyal, S.; Lin, Z.; Liu, J.; Guo, D.; Ott, M.; Zitnick, C. L.; Ma, J., et al. Biological structure and function emerge from scaling unsupervised learning

to 250 million protein sequences. *Proceedings of the National Academy of Sciences* **2021**, *118*, e2016239118.

(20) Hayes, T. et al. Simulating 500 million years of evolution with a language model. *bioRxiv* **2024**,

(21) Elnaggar, A.; Heinzinger, M.; Dallago, C.; Rehawi, G.; Wang, Y.; Jones, L.; Gibbs, T.; Feher, T.; Angerer, C.; Steinegger, M., et al. Prottrans: Toward understanding the language of life through self-supervised learning. *IEEE transactions on pattern analysis and machine intelligence* **2021**, *44*, 7112–7127.

(22) 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.

(23) Lv, L.; Lin, Z.; Li, H.; Liu, Y.; Cui, J.; Yu-Chian Chen, C.; Yuan, L.; Tian, Y. Prollama: A protein large language model for multi-task protein language processing. *arXiv e-prints* **2024**, arXiv–2402.

(24) Olsen, T. H.; Moal, I. H.; Deane, C. M. AbLang: an antibody language model for completing antibody sequences. *Bioinformatics Advances* **2022**, *2*, vbac046.

(25) Gao, K.; Wu, L.; Zhu, J.; Peng, T.; Xia, Y.; He, L.; Xie, S.; Qin, T.; Liu, H.; He, K., et al. Pre-training Antibody Language Models for Antigen-Specific Computational Antibody Design. Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining. 2023; pp 506–517.

(26) Shuai, R. W.; Ruffolo, J. A.; Gray, J. J. IgLM: Infilling language modeling for antibody sequence design. *Cell Systems* **2023**, *14*, 979–989.

(27) Kenlay, H.; Dreyer, F. A.; Kovaltsuk, A.; Miketa, D.; Pires, D.; Deane, C. M. Large scale paired antibody language models. *PLOS Computational Biology* **2024**, *20*, e1012646.

(28) Munsamy, G.; Lindner, S.; Lorenz, P.; Ferruz, N. ZymCTRL: a conditional language model for the controllable generation of artificial enzymes. NeurIPS machine learning in structural biology workshop. 2022.

(29) Hermosilla, P.; Ropinski, T. Contrastive representation learning for 3d protein structures. *arXiv preprint arXiv:2205.15675* **2022**,

(30) Xia, C.; Feng, S.-H.; Xia, Y.; Pan, X.; Shen, H.-B. Fast protein structure comparison through effective representation learning with contrastive graph neural networks. *PLoS computational biology* **2022**, *18*, e1009986.

(31) Zhou, B.; Lv, O.; Yi, K.; Xiong, X.; Tan, P.; Hong, L.; Wang, Y. G. Lightweight Equivariant Graph Representation Learning for Protein Engineering. **2022**,

(32) Le, T.; Noe, F.; Clevert, D.-A. Representation learning on biomolecular structures using equivariant graph attention. Learning on Graphs Conference. 2022; pp 30–1.

(33) Guo, Y.; Wu, J.; Ma, H.; Huang, J. Self-supervised pre-training for protein embeddings using tertiary structures. Proceedings of the AAAI conference on artificial intelligence. 2022; pp 6801–6809.

(34) Wang, L.; Liu, H.; Liu, Y.; Kurtin, J.; Ji, S. Learning Hierarchical Protein Representations via Complete 3D Graph Networks. The Eleventh International Conference on Learning Representations. 2023.

(35) Chen, C.; Chen, X.; Morehead, A.; Wu, T.; Cheng, J. 3D-equivariant graph neural networks for protein model quality assessment. *Bioinformatics* **2023**, *39*, btad030.

(36) Zhang, Z.; Xu, M.; Jamasb, A. R.; Chenthamarakshan, V.; Lozano, A.; Das, P.; Tang, J. Protein Representation Learning by Geometric Structure Pretraining. The Eleventh International Conference on Learning Representations. 2023.

(37) Gao, Z.; Jiang, C.; Zhang, J.; Jiang, X.; Li, L.; Zhao, P.; Yang, H.; Huang, Y.; Li, J. Hierarchical graph learning for protein–protein interaction. *Nature Communications* **2023**, *14*, 1093.

(38) Voitsitskyi, T.; Stratiichuk, R.; Koleiev, I.; Popryho, L.; Ostrovsky, Z.; Henitsoi, P.; Khropachov, I.; Vozniak, V.; Zhytar, R.; Nechepurenko, D., et al. 3DProtDTA: a deep learning model for drug-target affinity prediction based on residue-level protein graphs. *RSC advances* **2023**, *13*, 10261–10272.

(39) Wu, T.; Guo, Z.; Cheng, J. Atomic protein structure refinement using all-atom graph representations and SE (3)-equivariant graph transformer. *Bioinformatics* **2023**, *39*, btad298.

(40) jiale, Z.; Zhuang, W.; Song, J.; Li, Y.; Lu, S. Pre-Training Protein Bi-level Representation Through Span Mask Strategy On 3D Protein Chains. Forty-first International Conference on Machine Learning. 2024.

(41) Wang, Y.; Wu, S.; Duan, Y.; Huang, Y. A point cloud-based deep learning strategy for protein–ligand binding affinity prediction. *Briefings in bioinformatics* **2022**, *23*, bbab474.

(42) Nguyen, V. T. D.; Hy, T. S. Multimodal pretraining for unsupervised protein representation learning. *Biology Methods and Protocols* **2024**, *9*, bpae043.

(43) Gainza, P.; Sverrisson, F.; Monti, F.; Rodola, E.; Boscaini, D.; Bronstein, M. M.; Correia, B. E. Deciphering interaction fingerprints from protein molecular surfaces using geometric deep learning. *Nature Methods* **2020**, *17*, 184–192.

(44) Riahi, S.; Lee, J. H.; Sorenson, T.; Wei, S.; Jager, S.; Olfati-Saber, R.; Zhou, Y.; Park, A.; Wendt, M.; Minoux, H., et al. Surface ID: a geometry-aware system for protein molecular surface comparison. *Bioinformatics* **2023**, *39*, btad196.

(45) Mallet, V.; Attaiki, S.; Miao, Y.; Correia, B.; Ovsjanikov, M. AtomSurf: Surface Representation for Learning on Protein Structures. *arXiv preprint arXiv:2309.16519* **2023**,

(46) Randolph, N. Z.; Kuhlman, B. Invariant point message passing for protein side chain packing. *Proteins: Structure, Function, and Bioinformatics* **2024**,

(47) Marchand, A.; Buckley, S.; Schneuing, A.; Pacesa, M.; Elia, M.; Gainza, P.; Elizarova, E.; Neeser, R. M.; Lee, P.-W.; Reymond, L., et al. Targeting protein–ligand neosurfaces with a generalizable deep learning tool. *Nature* **2025**, 1–10.

(48) Sverrisson, F.; Feydy, J.; Correia, B. E.; Bronstein, M. M. Fast End-to-End Learning on Protein Surfaces. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). 2021; pp 15272–15281.

(49) Liu, Q.; Wang, P.-S.; Zhu, C.; Gaines, B. B.; Zhu, T.; Bi, J.; Song, M. OctSurf: Efficient hierarchical voxel-based molecular surface representation for protein-ligand affinity prediction. *Journal of Molecular Graphics and Modelling* **2021**, *105*, 107865.

(50) Srinivasan, B.; Ioannidis, V. N.; Adeshina, S.; Kakodkar, M.; Karypis, G.; Ribeiro, B. Conditional invariances for conformer invariant protein representations. **2022**,

(51) Sun, D.; Huang, H.; Li, Y.; Gong, X.; Ye, Q. DSR: Dynamical Surface Representation as Implicit Neural Networks for Protein. Thirty-seventh Conference on Neural Information Processing Systems. 2023.

(52) Satorras, V. G.; Hoogeboom, E.; Welling, M. E(n) Equivariant Graph Neural Networks. Proceedings of the 38th International Conference on Machine Learning. 2021; pp 9323–9332.

(53) Fuchs, F.; Worrall, D.; Fischer, V.; Welling, M. SE(3)-Transformers: 3D Roto-Translation Equivariant Attention Networks. Advances in Neural Information Processing Systems. 2020; pp 1970–1981.

(54) Anderson, B.; Hy, T. S.; Kondor, R. Cormorant: Covariant Molecular Neural Networks. Advances in Neural Information Processing Systems. 2019.

(55) Maron, H.; Ben-Hamu, H.; Shamir, N.; Lipman, Y. Invariant and Equivariant Graph Networks. International Conference on Learning Representations. 2019.

(56) Hy, T. S.; Trivedi, S.; Pan, H.; Anderson, B. M.; Kondor, R. Predicting molecular properties with covariant compositional networks. *The Journal of Chemical Physics* **2018**, *148*, 241745.

(57) Cheng, P.; Mao, C.; Tang, J.; Yang, S.; Cheng, Y.; Wang, W.; Gu, Q.; Han, W.; Chen, H.; Li, S., et al. Zero-shot prediction of mutation effects with multimodal deep representation learning guides protein engineering. *Cell Research* **2024**, *34*, 630–647.

(58) Wang, Z.; Combs, S. A.; Brand, R.; Calvo, M. R.; Xu, P.; Price, G.; Golovach, N.; Salawu, E. O.; Wise, C. J.; Ponnapalli, S. P., et al. Lm-gvp: an extensible sequence and structure informed deep learning framework for protein property prediction. *Scientific reports* **2022**, *12*, 6832.

(59) Gligorijević, V.; Renfrew, P. D.; Kosciolek, T.; Leman, J. K.; Berenberg, D.; Vatanen, T.; Chandler, C.; Taylor, B. C.; Fisk, I. M.; Vlamakis, H., et al. Structure-based protein function prediction using graph convolutional networks. *Nature communications* **2021**, *12*, 3168.

(60) Wu, F.; Wu, L.; Radev, D.; Xu, J.; Li, S. Z. Integration of pre-trained protein language models into geometric deep learning networks. *Communications Biology* **2023**, *6*, 876.

(61) Blaabjerg, L. M.; Jonsson, N.; Boomsma, W.; Stein, A.; Lindorff-Larsen, K. SSEmb: A joint embedding of protein sequence and structure enables robust variant effect predictions. *Nature Communications* **2024**, *15*, 9646.

(62) Ahmed, M.; Ali, S.; Jan, A.; Khan, I. U.; Patterson, M. Improved Graph-based Antibody-aware Epitope Prediction with Protein Language Model-based Embeddings. *bioRxiv* **2025**, 2025–02.

(63) Zheng, Z.; Deng, Y.; Xue, D.; Zhou, Y.; Ye, F.; Gu, Q. Structure-informed language models are protein designers. International conference on machine learning. 2023; pp 42317–42338.

(64) Yang, K. K.; Zanichelli, N.; Yeh, H. Masked inverse folding with sequence transfer for protein representation learning. *Protein Engineering, Design and Selection* **2023**, *36*, gzad015.

(65) Su, J.; Han, C.; Zhou, Y.; Shan, J.; Zhou, X.; Yuan, F. SaProt: Protein Language Modeling with Structure-aware Vocabulary. The Twelfth International Conference on Learning Representations. 2024.

(66) Sun, Y.; Shen, Y. Structure-informed protein language models are robust predictors for variant effects. *Human Genetics* **2024**, 1–17.

(67) Li, Z.; Cen, J.; Su, B.; Huang, W.; Xu, T.; Rong, Y.; Zhao, D. Large Language-Geometry Model: When LLM meets Equivariance. *arXiv preprint arXiv:2502.11149* **2025**,

(68) Lee, Y.; Yu, H.; Lee, J.; Kim, J. Pre-training Sequence, Structure, and Surface Features for Comprehensive Protein Representation Learning. The Twelfth International Conference on Learning Representations. 2023.

(69) Wang, D.; Pourmirzaei, M.; Abbas, U. L.; Zeng, S.; Manshour, N.; Esmaili, F.; Poudel, B.; Jiang, Y.; Shao, Q.; Chen, J.; Xu, D. S-PLM: Structure-Aware Protein Language Model via Contrastive Learning Between Sequence and Structure. *Advanced Science* **2025**, *12*, 2404212.

(70) Zhang, Z.; Xu, M.; Lozano, A. C.; Chenthamarakshan, V.; Das, P.; Tang, J. Pre-training protein encoder via siamese sequence-structure diffusion trajectory prediction. *Advances in Neural Information Processing Systems* **2024**, *36*.

(71) Barton, J.; Galson, J. D.; Leem, J. Enhancing antibody language models with structural information. *bioRxiv* **2024**, 2023–12.

(72) Ngo, N. K.; Hy, T. S. Multimodal protein representation learning and target-aware varia-tional auto-encoders for protein-binding ligand generation. *Machine Learning: Science and Technology* **2024**, *5*, 025021.

(73) Liu, H.; Jian, Y.; Zeng, C.; Zhao, Y. RNA-protein interaction prediction using network-guided deep learning. *Communications Biology* **2025**, *8*, 247.

(74) Lee, J.; Yoon, W.; Kim, S.; Kim, D.; Kim, S.; So, C. H.; Kang, J. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* **2020**, *36*, 1234–1240.

(75) Gu, Y.; Tinn, R.; Cheng, H.; Lucas, M.; Usuyama, N.; Liu, X.; Naumann, T.; Gao, J.; Poon, H. Domain-specific language model pretraining for biomedical natural language pro-cessing. *ACM Transactions on Computing for Healthcare (HEALTH)* **2021**, *3*, 1–23.

(76) Zhang, N.; Bi, Z.; Liang, X.; Cheng, S.; Hong, H.; Deng, S.; Zhang, Q.; Lian, J.; Chen, H. OntoProtein: Protein Pretraining With Gene Ontology Embedding. International Confer-ence on Learning Representations. 2022.

(77) Zhou, H.-Y.; Fu, Y.; Zhang, Z.; Cheng, B.; Yu, Y. Protein Representation Learning via Knowledge Enhanced Primary Structure Reasoning. The Eleventh International Conference on Learning Representations. 2023.

(78) Xu, M.; Yuan, X.; Miret, S.; Tang, J. Protst: Multi-modality learning of protein sequences and biomedical texts. International Conference on Machine Learning. 2023; pp 38749–38767.

(79) Kilgore, H. R.; Chinn, I.; Mikhael, P. G.; Mitnikov, I.; Van Dongen, C.; Zylberberg, G.; Afeyan, L.; Banani, S. F.; Wilson-Hawken, S.; Lee, T. I., et al. Protein codes promote selec-tive subcellular compartmentalization. *Science* **2025**, eadq2634.

(80) Hu, F.; Hu, Y.; Zhang, W.; Huang, H.; Pan, Y.; Yin, P. A multimodal protein representation framework for quantifying transferability across biochemical downstream tasks. *Advanced Science* **2023**, *10*, 2301223.

(81) Abdine, H.; Chatzianastasis, M.; Bouyioukos, C.; Vazirgiannis, M. Prot2text: Multimodal protein's function generation with gnns and transformers. Proceedings of the AAAI Conference on Artificial Intelligence. 2024; pp 10757–10765.

(82) 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.

(83) 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. *Journal of cheminformatics* **2022**, *14*, 14.

(84) Pei, Q.; Wu, L.; Zhu, J.; Xia, Y.; Xie, S.; Qin, T.; Liu, H.; Liu, T.-Y. Smt-dta: Improving drug-target affinity prediction with semi-supervised multi-task training. *arXiv preprint arXiv:2206.09818* **2022**,

(85) Gao, Z.; Tan, C.; Xia, J.; Li, S. Z. Co-supervised Pre-training of Pocket and Ligand. Joint European Conference on Machine Learning and Knowledge Discovery in Databases. 2023; pp 405–421.

(86) Nakata, S.; Mori, Y.; Tanaka, S. End-to-end protein–ligand complex structure generation with diffusion-based generative models. *BMC bioinformatics* **2023**, *24*, 233.

(87) Nguyen, V. T. D.; Nguyen, N. D.; Hy, T. S. ProteinReDiff: Complex-based ligand-binding proteins redesign by equivariant diffusion-based generative models. *Structural Dynamics* **2024**, *11*.

(88) Gao, B.; Qiang, B.; Tan, H.; Jia, Y.; Ren, M.; Lu, M.; Liu, J.; Ma, W.-Y.; Lan, Y. Drug-clip: Contrasive protein-molecule representation learning for virtual screening. *Advances in Neural Information Processing Systems* **2024**, *36*.

(89) Ferreira de Freitas, R.; Schapira, M. A systematic analysis of atomic protein–ligand interactions in the PDB. *Med. Chem. Commun.* **2017**, *8*, 1970–1981.

(90) Zheng, L.; Fan, J.; Mu, Y. Onionnet: a multiple-layer intermolecular-contact-based convolutional neural network for protein–ligand binding affinity prediction. *ACS omega* **2019**, *4*, 15956–15965.

(91) Moon, S.; Zhung, W.; Yang, S.; Lim, J.; Kim, W. Y. PIGNet: a physics-informed deep learning model toward generalized drug–target interaction predictions. *Chemical Science* **2022**, *13*, 3661–3673.

(92) Feng, S.; Li, M.; Jia, Y.; Ma, W.-Y.; Lan, Y. Protein-ligand binding representation learning from fine-grained interactions. The Twelfth International Conference on Learning Representations. 2024.

(93) Yuan, Y.; Chen, Q.; Mao, J.; Li, G.; Pan, X. DG-Affinity: predicting antigen-antibody affinity with language models from sequences. *BMC bioinformatics* **2023**, *24*, 430.

(94) Li, M.; Shi, Y.; Hu, S.; Hu, S.; Guo, P.; Wan, W.; Zhang, L. Y.; Pan, S.; Li, J.; Sun, L., et al. MVSF-AB: Accurate antibody-antigen binding affinity prediction via multi-view sequence feature learning. *Bioinformatics* **2024**, btae579.

(95) Bandara, N.; Premathilaka, D.; Chandanayake, S.; Hettiarachchi, S.; Varenthirarajah, V.; Munasinghe, A.; Madhawa, K.; Charles, S. Deep Geometric Framework to Predict Antibody-Antigen Binding Affinity. *bioRxiv* **2024**, 2024–06.

(96) Evans, R.; O'Neill, M.; Pritzel, A.; Antropova, N.; Senior, A.; Green, T.; Žídek, A.;

Bates, R.; Blackwell, S.; Yim, J., et al. Protein complex prediction with AlphaFold-Multimer. *biorxiv* **2021**, 2021–10.

(97) Giulini, M.; Schneider, C.; Cutting, D.; Desai, N.; Deane, C. M.; Bonvin, A. M. J. J. Towards the accurate modelling of antibody-antigen complexes from sequence using machine learning and information-driven docking. *Bioinformatics* **2024**, *40*, btae583.

(98) Ruffolo, J. A.; Chu, L.-S.; Mahajan, S. P.; Gray, J. J. Fast, accurate antibody structure prediction from deep learning on massive set of natural antibodies. *Nature communications* **2023**, *14*, 2389.

(99) Réau, M.; Renaud, N.; Xue, L. C.; Bonvin, A. M. DeepRank-GNN: a graph neural network framework to learn patterns in protein–protein interfaces. *Bioinformatics* **2023**, *39*, btac759.

(100) Yue, Y.; Li, S.; Cheng, Y.; Wang, L.; Hou, T.; Zhu, Z.; He, S. Integration of molecular coarse-grained model into geometric representation learning framework for protein-protein complex property prediction. *Nature Communications* **2024**, *15*, 9629.

(101) Wang, Z.; Brand, R.; Adolf-Bryfogle, J.; Grewal, J.; Qi, Y.; Combs, S. A.; Golovach, N.; Alford, R.; Rangwala, H.; Clark, P. M. EGGNet, a generalizable geometric deep learning framework for protein complex pose scoring. *ACS omega* **2024**, *9*, 7471–7479.

(102) Boutet, E.; Lieberherr, D.; Tognolli, M.; Schneider, M.; Bairoch, A. *Plant bioinformatics: methods and protocols*; Springer, 2007; pp 89–112.

(103) Möller, S.; Leser, U.; Fleischmann, W.; Apweiler, R. EDITtoTrEMBL: a distributed approach to high-quality automated protein sequence annotation. *Bioinformatics* **1999**, *15*, 219–227.

(104) Suzek, B. E.; Huang, H.; McGarvey, P.; Mazumder, R.; Wu, C. H. UniRef: comprehensive and non-redundant UniProt reference clusters. *Bioinformatics* **2007**, *23*, 1282–1288.

(105) Mistry, J.; Chuguransky, S.; Williams, L.; Qureshi, M.; Salazar, G. A.; Sonnhammer, E. L.; Tosatto, S. C.; Paladin, L.; Raj, S.; Richardson, L. J., et al. Pfam: The protein families database in 2021. *Nucleic acids research* **2021**, *49*, D412–D419.

(106) Richardson, L.; Allen, B.; Baldi, G.; Beracochea, M.; Bileschi, M. L.; Burdett, T.; Burgin, J.; Caballero-Pérez, J.; Cochrane, G.; Colwell, L. J., et al. MGnify: the microbiome sequence data analysis resource in 2023. *Nucleic Acids Research* **2023**, *51*, D753–D759.

(107) 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 research* **2000**, *28*, 235–242.

(108) Varadi, M.; Bertoni, D.; Magana, P.; Paramval, U.; Pidruchna, I.; Radhakrishnan, M.; Tsenkov, M.; Nair, S.; Mirdita, M.; Yeo, J., et al. AlphaFold Protein Structure Database in 2024: providing structure coverage for over 214 million protein sequences. *Nucleic acids research* **2024**, *52*, D368–D375.

(109) Aleksander, S. A.; Balhoff, J.; Carbon, S.; Cherry, J. M.; Drabkin, H. J.; Ebert, D.; Feuermann, M.; Gaudet, P.; Harris, N. L., et al. The gene ontology knowledgebase in 2023. *Genetics* **2023**, *224*, iyad031.

(110) Rao, R.; Bhattacharya, N.; Thomas, N.; Duan, Y.; Chen, P.; Canny, J.; Abbeel, P.; Song, Y. Evaluating protein transfer learning with TAPE. *Advances in neural information processing systems* **2019**, *32*.

(111) Liu, Z.; Su, M.; Han, L.; Liu, J.; Yang, Q.; Li, Y.; Wang, R. Forging the basis for developing protein–ligand interaction scoring functions. *Accounts of chemical research* **2017**, *50*, 302–309.

(112) Li, J.; Guan, X.; Zhang, O.; Sun, K.; Wang, Y.; Bagni, D.; Head-Gordon, T. Leak Proof PDBBind: A Reorganized Dataset of Protein-Ligand Complexes for More Generalizable Binding Affinity Prediction. 2024; https://arxiv.org/abs/2308.09639.

(113) 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. *Nature biotechnology* **2011**, *29*, 1046–1051.

(114) Tang, J.; Szwajda, A.; Shakyawar, S.; Xu, T.; Hintsanen, P.; Wennerberg, K.; Aittokallio, T. Making Sense of Large-Scale Kinase Inhibitor Bioactivity Data Sets: A Comparative and Integrative Analysis. *Journal of Chemical Information and Modeling* **2014**, *54*, 735–743, PMID: 24521231.

(115) Durairaj, J. et al. PLINDER: The protein-ligand interactions dataset and evaluation resource. *bioRxiv* **2024**,

(116) Kovtun, D.; Akdel, M.; Goncearenco, A.; Zhou, G.; Holt, G.; Baugher, D.; Lin, D.; Adeshina, Y.; Castiglione, T.; Wang, X., et al. PINDER: The protein interaction dataset and evaluation resource. *bioRxiv* **2024**, 2024–07.

(117) Yu, T.; Cui, H.; Li, J. C.; Luo, Y.; Jiang, G.; Zhao, H. Enzyme function prediction using contrastive learning. *Science* **2023**, *379*, 1358–1363.

(118) Buton, N.; Coste, F.; Le Cunff, Y. Predicting enzymatic function of protein sequences with attention. *Bioinformatics* **2023**, *39*, btad620.

(119) Yuan, Y.; Chen, S.; Hu, R.; Wang, X. MutualDTA: An Interpretable Drug–Target Affinity Prediction Model Leveraging Pretrained Models and Mutual Attention. *Journal of Chemical Information and Modeling* **2025**,

(120) Ha, C. N.; Pham, P.; Hy, T. S. LANTERN: Leveraging Large Language Models and Transformers for Enhanced Molecular Interactions. *bioRxiv* **2025**, 2025–02.

(121) Luo, D.; Liu, D.; Qu, X.; Dong, L.; Wang, B. Enhancing generalizability in protein–ligand binding affinity prediction with multimodal contrastive learning. *Journal of Chemical Information and Modeling* **2024**, *64*, 1892–1906.

(122) Jaskolski, M.; Dauter, Z.; Wlodawer, A. A brief history of macromolecular crystallography, illustrated by a family tree and its Nobel fruits. *The FEBS Journal* **2014**, *281*, 3985–4009.

(123) chen Bai, X.; McMullan, G.; Scheres, S. H. How cryo-EM is revolutionizing structural biology. *Trends in Biochemical Sciences* **2015**, *40*, 49–57.

(124) Vénien-Bryan, C.; Li, Z.; Vuillard, L.; Boutin, J. A. Cryo-electron microscopy and X-ray crystallography: complementary approaches to structural biology and drug discovery. *Structural Biology and Crystallization Communications* **2017**, *73*, 174–183.

(125) Kryshtafovych, A.; Schwede, T.; Topf, M.; Fidelis, K.; Moult, J. Critical assessment of methods of protein structure prediction (CASP)—Round XIV. *Proteins: Structure, Function, and Bioinformatics* **2021**, *89*, 1607–1617.

(126) Robin, X.; Haas, J.; Gumienny, R.; Smolinski, A.; Tauriello, G.; Schwede, T. Continuous Automated Model EvaluatiOn (CAMEO)—Perspectives on the future of fully automated evaluation of structure prediction methods. *Proteins: Structure, Function, and Bioinformatics* **2021**, *89*, 1977–1986.

(127) Baek, M.; DiMaio, F.; Anishchenko, I.; Dauparas, J.; Ovchinnikov, S.; Lee, G. R.; Wang, J.; Cong, Q.; Kinch, L. N.; Schaeffer, R. D., et al. Accurate prediction of protein structures and interactions using a three-track neural network. *Science* **2021**, *373*, 871–876.

(128) Ahdritz, G.; Bouatta, N.; Floristean, C.; Kadyan, S.; Xia, Q.; Gerecke, W.; O'Donnell, T. J.; Berenberg, D.; Fisk, I.; Zanichelli, N., et al. OpenFold: Retraining AlphaFold2 yields new insights into its learning mechanisms and capacity for generalization. *Nature Methods* **2024**, *21*, 1514–1524.

(129) Wu, R.; Ding, F.; Wang, R.; Shen, R.; Zhang, X.; Luo, S.; Su, C.; Wu, Z.; Xie, Q.; Berger, B., et al. High-resolution de novo structure prediction from primary sequence. *BioRxiv* **2022**, 2022–07.

(130) Fang, X.; Wang, F.; Liu, L.; He, J.; Lin, D.; Xiang, Y.; Zhu, K.; Zhang, X.; Wu, H.; Li, H., et al. A method for multiple-sequence-alignment-free protein structure prediction using a protein language model. *Nature Machine Intelligence* **2023**, *5*, 1087–1096.

(131) Lewis, S.; Hempel, T.; Jiménez-Luna, J.; Gastegger, M.; Xie, Y.; Foong, A. Y.; Satorras, V. G.; Abdin, O.; Veeling, B. S.; Zaporozhets, I., et al. Scalable emulation of protein equilibrium ensembles with generative deep learning. *bioRxiv* **2024**, 2024–12.

(132) Arnold, F. H. Design by Directed Evolution. *Accounts of Chemical Research* **1998**, *31*, 125–131.

(133) Tracewell, C. A.; Arnold, F. H. Directed enzyme evolution: climbing fitness peaks one amino acid at a time. *Current Opinion in Chemical Biology* **2009**, *13*, 3–9, Biocatalysis and Biotransformation/Bioinorganic Chemistry.

(134) Romero, P. A.; Krause, A.; Arnold, F. H. Navigating the protein fitness landscape with Gaussian processes. *Proceedings of the National Academy of Sciences* **2013**, *110*, E193–E201.

(135) Sarkisyan, K. S. et al. Local fitness landscape of the green fluorescent protein. *Nature* **2016**, *533*, 397–401.

(136) Bryant, D. H.; Bashir, A.; Sinai, S.; Jain, N. K.; Ogden, P. J.; Riley, P. F.; Church, G. M.; Colwell, L. J.; Kelsic, E. D. Deep diversification of an AAV capsid protein by machine learning. *Nature Biotechnology* **2021**, *39*, 691–696.

(137) Firnberg, E.; Labonte, J. W.; Gray, J. J.; Ostermeier, M. A Comprehensive, High-Resolution Map of a Gene's Fitness Landscape. *Molecular Biology and Evolution* **2014**, *31*, 1581–1592.

(138) Starita, L. M.; Pruneda, J. N.; Lo, R. S.; Fowler, D. M.; Kim, H. J.; Hiatt, J. B.; Shendure, J.; Brzovic, P. S.; Fields, S.; Klevit, R. E. Activity-enhancing mutations in an E3 ubiquitin

ligase identified by high-throughput mutagenesis. *Proceedings of the National Academy of Sciences* **2013**, *110*, E1263–E1272.

(139) Wrenbeck, E. E.; Azouz, L. R.; Whitehead, T. A. Single-mutation fitness landscapes for an enzyme on multiple substrates reveal specificity is globally encoded. *Nature Communications* **2017**, *8*, 15695.

(140) Klesmith, J. R.; Bacik, J.-P.; Michalczyk, R.; Whitehead, T. A. Comprehensive Sequence-Flux Mapping of a Levoglucosan Utilization Pathway in E. coli. *ACS Synthetic Biology* **2015**, *4*, 1235–1243.

(141) Melamed, D.; Young, D. L.; Gamble, C. E.; Miller, C. R.; Fields, S. Deep mutational scanning of an RRM domain of the Saccharomyces cerevisiae poly(A)-binding protein. *RNA* **2013**, *19*, 1537–1551.

(142) Weile, J. et al. A framework for exhaustively mapping functional missense variants. *Molecular Systems Biology* **2017**, *13*, 957.

(143) Emami, P.; Perreault, A.; Law, J.; Biagioni, D.; John, P. S. Plug & play directed evolution of proteins with gradient-based discrete MCMC. *Machine Learning: Science and Technology* **2023**, *4*, 025014.

(144) Tran, T. V.; Hy, T. S. Protein design by directed evolution guided by large language models. *IEEE Transactions on Evolutionary Computation* **2024**,

(145) Castro, E.; Godavarthi, A.; Rubinfien, J.; Givechian, K.; Bhaskar, D.; Krishnaswamy, S. Transformer-based protein generation with regularized latent space optimization. *Nature Machine Intelligence* **2022**, *4*, 840–851.

(146) Tran, T.; Ngo, N. K.; Nguyen, V. T. D.; Hy, T.-S. LatentDE: Latent-based Directed Evolution for Protein Sequence Design. *Machine Learning: Science and Technology* **2025**,

(147) Kirjner, A.; Yim, J.; Samusevich, R.; Bracha, S.; Jaakkola, T.; Barzilay, R.; Fiete, I. Improving protein optimization with smoothed fitness landscapes. *arXiv preprint arXiv:2307.00494* **2023**,

(148) Tran, T. V.; Ngo, N. K.; Nguyen, V. A.; Hy, T. S. GROOT: Effective Design of Biological Sequences with Limited Experimental Data. *arXiv preprint arXiv:2411.11265* **2024**,

(149) Dunbar Jr, J. B.; Smith, R. D.; Damm-Ganamet, K. L.; Ahmed, A.; Esposito, E. X.; Delproposto, J.; Chinnaswamy, K.; Kang, Y.-N.; Kubish, G.; Gestwicki, J. E., et al. CSAR data set release 2012: ligands, affinities, complexes, and docking decoys. *Journal of chemical information and modeling* **2013**, *53*, 1842–1852.

(150) Polizzi, N. F.; DeGrado, W. F. A defined structural unit enables de novo design of small-molecule–binding proteins. *Science* **2020**, *369*, 1227–1233.

(151) Stark, H.; Jing, B.; Barzilay, R.; Jaakkola, T. Harmonic Prior Self-conditioned Flow Matching for Multi-Ligand Docking and Binding Site Design. NeurIPS 2023 AI for Science Workshop. 2023.

(152) Watson, J. L.; Juergens, D.; Bennett, N. R.; Trippe, B. L.; Yim, J.; Eisenach, H. E.; Ahern, W.; Borst, A. J.; Ragotte, R. J.; Milles, L. F., et al. De novo design of protein structure and function with RFdiffusion. *Nature* **2023**, *620*, 1089–1100.

(153) Dauparas, J.; Lee, G. R.; Pecoraro, R.; An, L.; Anishchenko, I.; Glasscock, C.; Baker, D. Atomic context-conditioned protein sequence design using LigandMPNN. *Biorxiv* **2023**, 2023–12.

(154) Hsu, C.; Verkuil, R.; Liu, J.; Lin, Z.; Hie, B.; Sercu, T.; Lerer, A.; Rives, A. Learning inverse folding from millions of predicted structures. International conference on machine learning. 2022; pp 8946–8970.

(155) Gao, Z.; Tan, C.; Li, S. Z. PiFold: Toward effective and efficient protein inverse folding. The Eleventh International Conference on Learning Representations. 2023.

(156) Meller, A.; Ward, M.; Borowsky, J.; Kshirsagar, M.; Lotthammer, J. M.; Oviedo, F.; Ferres, J. L.; Bowman, G. R. Predicting locations of cryptic pockets from single protein structures using the PocketMiner graph neural network. *Nature Communications* **2023**, *14*, 1177.

(157) Gagliardi, L.; Rocchia, W. SiteFerret: Beyond simple pocket identification in proteins. *J. Chem. Theory Comput.* **2023**, *19*, 5242–5259.

(158) Olsen, T. H.; Boyles, F.; Deane, C. M. Observed Antibody Space: A diverse database of cleaned, annotated, and translated unpaired and paired antibody sequences. *Protein Science* **2022**, *31*, 141–146.

(159) Dunbar, J.; Krawczyk, K.; Leem, J.; Baker, T.; Fuchs, A.; Georges, G.; Shi, J.; Deane, C. M. SAbDab: the structural antibody database. *Nucleic acids research* **2014**, *42*, D1140–D1146.

(160) Sirin, S.; Apgar, J. R.; Bennett, E. M.; Keating, A. E. AB-bind: antibody binding mutational database for computational affinity predictions. *Protein Science* **2016**, *25*, 393–409.

(161) Eguchi, R. R.; Choe, C. A.; Huang, P.-S. Ig-VAE: Generative modeling of protein structure by direct 3D coordinate generation. *PLoS computational biology* **2022**, *18*, e1010271.

(162) Shan, S.; Luo, S.; Yang, Z.; Hong, J.; Su, Y.; Ding, F.; Fu, L.; Li, C.; Chen, P.; Ma, J., et al. Deep learning guided optimization of human antibody against SARS-CoV-2 variants with broad neutralization. *Proceedings of the National Academy of Sciences* **2022**, *119*, e2122954119.

(163) Cohen, T.; Schneidman-Duhovny, D. Epitope-specific antibody design using diffusion models on the latent space of ESM embeddings. NeurIPS 2023 Generative AI and Biology (GenBio) Workshop. 2023.

(164) Wu, L.; Liu, Y.; Lin, H.; Huang, Y.; Zhao, G.; Gao, Z.; Li, S. Z. A Simple yet Effective $\Delta\Delta G$ Predictor is An Unsupervised Antibody Optimizer and Explainer. The Thirteenth International Conference on Learning Representations. 2025.

(165) Lin, J. Y.-Y.; Hofmann, J. L.; Leaver-Fay, A.; Liang, W.-C.; Vasilaki, S.; Lee, E.; Pinheiro, P. O.; Tagasovska, N.; Kiefer, J. R.; Wu, Y., et al. DyAb: sequence-based antibody design and property prediction in a low-data regime. *bioRxiv* **2025**, 2025–01.

(166) Jin, W.; Wohlwend, J.; Barzilay, R.; Jaakkola, T. S. Iterative Refinement Graph Neural Network for Antibody Sequence-Structure Co-design. International Conference on Learning Representations. 2022.

(167) Luo, S.; Su, Y.; Peng, X.; Wang, S.; Peng, J.; Ma, J. Antigen-specific antibody design and optimization with diffusion-based generative models for protein structures. *Advances in Neural Information Processing Systems* **2022**, *35*, 9754–9767.

(168) Villegas-Morcillo, A.; Weber, J.; Reinders, M. Guiding diffusion models for antibody sequence and structure co-design with developability properties. NeurIPS 2023 Generative AI and Biology (GenBio) Workshop. 2023.

(169) Martinkus, K.; Ludwiczak, J.; LIANG, W.-C.; Lafrance-Vanasse, J.; Hotzel, I.; Rajpal, A.; Wu, Y.; Cho, K.; Bonneau, R.; Gligorijevic, V.; Loukas, A. AbDiffuser: full-atom generation of in-vitro functioning antibodies. Thirty-seventh Conference on Neural Information Processing Systems. 2023.

(170) Zhou, X.; Xue, D.; Chen, R.; Zheng, Z.; Wang, L.; Gu, Q. Antigen-specific antibody design via direct energy-based preference optimization. *Advances in Neural Information Processing Systems* **2024**, *37*, 120861–120891.

(171) Yang, J.; Li, F.-Z.; Arnold, F. H. Opportunities and Challenges for Machine Learning-Assisted Enzyme Engineering. *ACS Central Science* **2024**, *10*, 226–241.

(172) Chang, A.; Jeske, L.; Ulbrich, S.; Hofmann, J.; Koblitz, J.; Schomburg, I.; Neumann-Schaal, M.; Jahn, D.; Schomburg, D. BRENDA, the ELIXIR core data resource in 2021: new developments and updates. *Nucleic acids research* **2021**, *49*, D498–D508.

(173) Yang, J.; Bhatnagar, A.; Ruffolo, J. A.; Madani, A. Conditional enzyme generation using protein language models with adapters. *arXiv preprint arXiv:2410.03634* **2024**,

(174) Hossack, E. J.; Hardy, F. J.; Green, A. P. Building enzymes through design and evolution. *ACS Catalysis* **2023**, *13*, 12436–12444.

(175) Mikhael, P. G.; Chinn, I.; Barzilay, R. Clipzyme: Reaction-conditioned virtual screening of enzymes. *arXiv preprint arXiv:2402.06748* **2024**,

(176) Hua, C.; Zhong, B.; Luan, S.; Hong, L.; Wolf, G.; Precup, D.; Zheng, S. Reactzyme: A benchmark for enzyme-reaction prediction. *Advances in Neural Information Processing Systems* **2025**, *37*, 26415–26442.

(177) Hua, C.; Lu, J.; Liu, Y.; Zhang, O.; Tang, J.; Ying, R.; Jin, W.; Wolf, G.; Precup, D.; Zheng, S. Reaction-conditioned De Novo Enzyme Design with GENzyme. *arXiv preprint arXiv:2411.16694* **2024**,

(178) Francoeur, P. G.; Masuda, T.; Sunseri, J.; Jia, A.; Iovanisci, R. B.; Snyder, I.; Koes, D. R. Three-dimensional convolutional neural networks and a cross-docked data set for structure-based drug design. *Journal of chemical information and modeling* **2020**, *60*, 4200–4215.

(179) Benson, M. L.; Smith, R. D.; Khazanov, N. A.; Dimcheff, B.; Beaver, J.; Dresslar, P.; Nerothin, J.; Carlson, H. A. Binding MOAD, a high-quality protein–ligand database. *Nucleic acids research* **2007**, *36*, D674–D678.

(180) Peng, X.; Luo, S.; Guan, J.; Xie, Q.; Peng, J.; Ma, J. Pocket2mol: Efficient molecular sampling based on 3d protein pockets. International Conference on Machine Learning. 2022; pp 17644–17655.

(181) Zhang, O.; Zhang, J.; Jin, J.; Zhang, X.; Hu, R.; Shen, C.; Cao, H.; Du, H.; Kang, Y.; Deng, Y., et al. ResGen is a pocket-aware 3D molecular generation model based on parallel multiscale modelling. *Nature Machine Intelligence* **2023**, *5*, 1020–1030.

(182) Jiang, Y.; Zhang, G.; You, J.; Zhang, H.; Yao, R.; Xie, H.; Zhang, L.; Xia, Z.; Dai, M.; Wu, Y., et al. Pocketflow is a data-and-knowledge-driven structure-based molecular generative model. *Nature Machine Intelligence* **2024**, *6*, 326–337.

(183) Zhung, W.; Kim, H.; Kim, W. Y. 3D molecular generative framework for interaction-guided drug design. *Nature Communications* **2024**, *15*, 2688.

(184) Schneuing, A.; Harris, C.; Du, Y.; Didi, K.; Jamasb, A.; Igashov, I.; Du, W.; Gomes, C.; Blundell, T. L.; Lio, P., et al. Structure-based drug design with equivariant diffusion models. *Nature Computational Science* **2024**, *4*, 899–909.

(185) Chen, L.; Fan, Z.; Chang, J.; Yang, R.; Hou, H.; Guo, H.; Zhang, Y.; Yang, T.; Zhou, C.; Sui, Q., et al. Sequence-based drug design as a concept in computational drug design. *Nature communications* **2023**, *14*, 4217.

(186) Creanza, T. M.; Alberga, D.; Patruno, C.; Mangiatordi, G. F.; Ancona, N. Transformer Decoder Learns from a Pretrained Protein Language Model to Generate Ligands with High Affinity. *Journal of Chemical Information and Modeling* **2025**,

(187) 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. *Journal of medicinal chemistry* **2012**, *55*, 6582–6594.

(188) Tran-Nguyen, V.-K.; Jacquemard, C.; Rognan, D. LIT-PCBA: an unbiased data set for machine learning and virtual screening. *Journal of chemical information and modeling* **2020**, *60*, 4263–4273.

(189) Chandak, P.; Huang, K.; Zitnik, M. Building a knowledge graph to enable precision medicine. *Scientific Data* **2023**, *10*, 67.

(190) Huang, K.; Fu, T.; Gao, W.; Zhao, Y.; Roohani, Y. H.; Leskovec, J.; Coley, C. W.; Xiao, C.; Sun, J.; Zitnik, M. Therapeutics Data Commons: Machine Learning Datasets and Tasks for Drug Discovery and Development. Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1). 2021.

(191) Zagidullin, B.; Aldahdooh, J.; Zheng, S.; Wang, W.; Wang, Y.; Saad, J.; Malyutina, A.; Jafari, M.; Tanoli, Z.; Pessia, A., et al. DrugComb: an integrative cancer drug combination data portal. *Nucleic acids research* **2019**, *47*, W43–W51.

(192) Zitnik, M.; Sosič, R.; Maheshwari, S.; Leskovec, J. BioSNAP Datasets: Stanford Biomedical Network Dataset Collection. 2018.

(193) Trott, O.; Olson, A. J. AutoDock Vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *Journal of computational chemistry* **2010**, *31*, 455–461.

(194) Krivák, R.; Hoksza, D. P2Rank: machine learning based tool for rapid and accurate prediction of ligand binding sites from protein structure. *Journal of cheminformatics* **2018**, *10*, 1–12.

(195) Stärk, H.; Ganea, O.; Pattanaik, L.; Barzilay, R.; Jaakkola, T. Equibind: Geometric deep learning for drug binding structure prediction. International conference on machine learning. 2022; pp 20503–20521.

(196) Lu, W.; Wu, Q.; Zhang, J.; Rao, J.; Li, C.; Zheng, S. Tankbind: Trigonometry-aware neural networks for drug-protein binding structure prediction. *Advances in neural information processing systems* **2022**, *35*, 7236–7249.

(197) Corso, G.; Stärk, H.; Jing, B.; Barzilay, R.; Jaakkola, T. S. DiffDock: Diffusion Steps, Twists, and Turns for Molecular Docking. The Eleventh International Conference on Learning Representations. 2023.

(198) Liu, L.; Zhang, S.; He, D.; Ye, X.; Zhou, J.; Zhang, X.; Jiang, Y.; Diao, W.; Yin, H.; Chai, H., et al. Pre-training on large-scale generated docking conformations with helixdock to unlock the potential of protein-ligand structure prediction models. *arXiv preprint arXiv:2310.13913* **2023**,

(199) Bender, A.; Jenkins, J. L.; Scheiber, J.; Sukuru, S. C. K.; Glick, M.; Davies, J. W. How similar are similarity searching methods? A principal component analysis of molecular descriptor space. *Journal of chemical information and modeling* **2009**, *49*, 108–119.

(200) Singh, R.; Sledzieski, S.; Bryson, B.; Cowen, L.; Berger, B. Contrastive learning in protein language space predicts interactions between drugs and protein targets. *Proceedings of the National Academy of Sciences* **2023**, *120*, e2220778120.

(201) Gao, B.; Qiang, B.; Tan, H.; Jia, Y.; Ren, M.; Lu, M.; Liu, J.; Ma, W.-Y.; Lan, Y. DrugCLIP: Contrasive Protein-Molecule Representation Learning for Virtual Screening. Thirty-seventh Conference on Neural Information Processing Systems. 2023.

(202) McNutt, A. T.; Adduri, A. K.; Ellington, C. N.; Dayao, M. T.; Xing, E. P.; Mohimani, H.; Koes, D. R. SPRINT Enables Interpretable and Ultra-Fast Virtual Screening against Thousands of Proteomes. **2024**,

(203) Daza, D.; Alivanistos, D.; Mitra, P.; Pijnenburg, T.; Cochez, M.; Groth, P. BioBLP: a modular framework for learning on multimodal biomedical knowledge graphs. *Journal of Biomedical Semantics* **2023**, *14*, 20.

(204) Djeddi, W. E.; Hermi, K.; Ben Yahia, S.; Diallo, G. Advancing drug–target interaction prediction: a comprehensive graph-based approach integrating knowledge graph embedding and ProtBert pretraining. *BMC bioinformatics* **2023**, *24*, 488.

(205) Dang, T.; Nguyen, V. T. D.; Le, M. T.; Hy, T.-S. Multimodal Contrastive Representation Learning in Augmented Biomedical Knowledge Graphs. *arXiv preprint arXiv:2501.01644* **2025**,

(206) Dong, W.; Yang, Q.; Wang, J.; Xu, L.; Li, X.; Luo, G.; Gao, X. Multi-modality attribute learning-based method for drug–protein interaction prediction based on deep neural network. *Briefings in bioinformatics* **2023**, *24*, bbad161.

(207) Nguyen, E.; Poli, M.; Durrant, M. G.; Kang, B.; Katrekar, D.; Li, D. B.; Bartie, L. J.; Thomas, A. W.; King, S. H.; Brixi, G., et al. Sequence modeling and design from molecular to genome scale with Evo. *Science* **2024**, *386*, eado9336.

(208) Brixi, G.; Durrant, M. G.; Ku, J.; Poli, M.; Brockman, G.; Chang, D.; Gonzalez, G. A.; King, S. H.; Li, D. B.; Merchant, A. T., et al. Genome modeling and design across all domains of life with Evo 2. *bioRxiv* **2025**, 2025–02.

(209) Consens, M. E.; Dufault, C.; Wainberg, M.; Forster, D.; Karimzadeh, M.; Goodarzi, H.; Theis, F. J.; Moses, A.; Wang, B. Transformers and genome language models. *Nature Machine Intelligence* **2025**,

(210) Sanh, V.; Debut, L.; Chaumond, J.; Wolf, T. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* **2019**,

(211) Jiao, X.; Yin, Y.; Shang, L.; Jiang, X.; Chen, X.; Li, L.; Wang, F.; Liu, Q. Tinybert: Distilling bert for natural language understanding. *arXiv preprint arXiv:1909.10351* **2019**,

(212) Wu, K.; Zhang, J.; Peng, H.; Liu, M.; Xiao, B.; Fu, J.; Yuan, L. Tinyvit: Fast pretraining distillation for small vision transformers. European conference on computer vision. 2022; pp 68–85.

(213) Chen, X.; Cao, Q.; Zhong, Y.; Zhang, J.; Gao, S.; Tao, D. Dearkd: data-efficient early

knowledge distillation for vision transformers. Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2022; pp 12052–12062.

(214) Tan, Y.; Wang, R.; Wu, B.; Hong, L.; Zhou, B. Retrieval-enhanced mutation mastery: Augmenting zero-shot prediction of protein language model. *arXiv preprint arXiv:2410.21127* **2024**,

(215) Gong, X.; Liu, Q.; He, J.; Guo, Y.; Wang, G. MultiGranDTI: an explainable multi-granularity representation framework for drug-target interaction prediction. *Applied Intelligence* **2025**, *55*, 1–19.

(216) Medina-Ortiz, D.; Khalifeh, A.; Anvari-Kazemabad, H.; Davari, M. D. Interpretable and explainable predictive machine learning models for data-driven protein engineering. *Biotechnology Advances* **2024**, 108495.

(217) Ma, Z.; Fan, C.; Wang, Z.; Chen, Z.; Lin, X.; Li, Y.; Feng, S.; Zhang, J.; Cao, Z.; Gao, Y. Q. ProtTeX: Structure-In-Context Reasoning and Editing of Proteins with Large Language Models. 2025; `https://arxiv.org/abs/2503.08179`.

(218) Wang, Z.; Ma, Z.; Cao, Z.; Zhou, C.; Zhang, J.; Gao, Y. Prot2Chat: Protein LLM with Early Fusion of Sequence and Structure. 2025; `https://arxiv.org/abs/2502.06846`.