

# Computational Protein Science in the Era of Large Language Models (LLMs)

Wenqi Fan, Yi Zhou, Shijie Wang, Yuyao Yan, Hui Liu, Qian Zhao, Le Song, and Qing Li

**Abstract**—Proteins are macromolecules that play essential roles in almost all essential life activities, such as immunity, digestion, disease regulation, etc. Considering the significance of proteins, computational protein science has always been a critical field of scientific research, dedicated to revealing knowledge and developing applications within the protein sequence-structure-function paradigm. In the last few decades, Artificial Intelligence (AI) has made a significant impact in computational protein science, leading to notable, even Nobel-Prize-level successes in various specific protein modeling tasks. However, those previous AI models still meet limitations, such as the difficulty in comprehending the grammar and semantics contained in protein sequences, and the inability to generalize across a wide range of protein modeling tasks. Recently, Large Language Models (LLMs) have emerged as a milestone in AI advances due to their remarkable language processing capability and unprecedented generalization capability. They are capable of promoting comprehensive progress in fields, rather than merely solving individual tasks. As a result, researchers have actively introduced powerful LLM techniques in promoting computational protein science, developing protein Language Models (pLMs) that skillfully grasp the foundational knowledge of proteins and can be effectively generalized to solve a diversity of sequence-structure-function reasoning problems. While witnessing prosperous developments, it's necessary to present a systematic overview of computational protein science empowered by LLM techniques. First, we summarize existing pLMs into categories based on their mastered protein knowledge, i.e., underlying sequence patterns, explicit structural and functional information, and external scientific languages. Second, we introduce the utilization and adaptation of pLMs, highlighting their remarkable achievements in promoting protein structure prediction, protein function prediction, and protein design studies. Then, we describe the practical application of pLMs in antibody design, enzyme design, and drug discovery. Finally, we specifically discuss the promising future directions in this fast-growing field.

**Index Terms**—Protein Language Models, Protein Structure Prediction, Protein Function Prediction, Protein Design, and Large Language Models (LLMs).

## 1 INTRODUCTION

As the most foundational building blocks of life, proteins play essential roles in almost all biological cellular processes [1, 2], such as metabolism, signal transduction, immune responses, etc. After long research, as illustrated in Figure 1, people have reached a limited understanding of the nature of proteins: Primarily, proteins adhere to the **sequence-structure-function** paradigm [3, 4] — the *amino acid (AA)<sup>1</sup> sequence of a protein indicates its three-dimensional structure, which in turn determines its function*. Moreover, proteins are shaped by the forces of **evolution** — natural selection reserves protein sequences capable of folding into stable structures and fulfilling proper functions, while eliminating those that do not [5]. Therefore,

the protein sequence is widely acknowledged as the protein language [6, 2], where the underlying arrangement patterns of AAs resemble "grammar" and the encoded structural and functional information mirrors "semantics." In this context, in the continued progression of scientific exploration, there are greater challenges in deciphering the protein language and applying rules of information flow among the protein sequence-structure-function. Computational protein science has emerged as a vitally important research field.

With the ability to extract patterns and fit mappings from data, Artificial Intelligence (AI) techniques have been widely adopted in computational protein science, which has even driven groundbreaking results at the Nobel Prize level, with David Baker awarded for "*computational protein design*" and Demis Hassabis & John Jumper jointly awarded for "*protein structure prediction*". Technically, the "AI for Protein" studies propose a variety of networks to accomplish different protein modeling tasks. For instance, UniRep [7] trains a mLSTM model on unlabeled AA sequences to distill the fundamental features of a protein into a statistical representation. AlphaFold2 [8] and RoseTTAFold [9] achieve a breakthrough in the accurate prediction of protein structures by exploiting the evolutionary information within multiple homologous sequences. DeepGOPuls [10] combines heuristic sequence similarity and a CNN model to predict protein functional annotations. Variational Autoencoder (VAE) is employed for protein sequence generation [11] and backbone structure generation [12], which are important links in computational protein design. Although these AI-empowered

- W. Fan is with the Department of Computing and Department of Management and Marketing, The Hong Kong Polytechnic University. E-mail: wenqifan03@gmail.com.
- Y. Zhou, S. Wang, and Q. Li are with the Department of Computing, The Hong Kong Polytechnic University. Email: {echo-yi.zhou, shijie.wang}@connect.polyu.hk, csqli@comp.polyu.edu.hk.
- Y. Yan is with the Department of Health Technology and Informatics, The Hong Kong Polytechnic University. E-mail: yanyuyao69@gmail.com.
- Hui Liu is with Michigan State University. Email: liuhui7@msu.edu.
- Qian Zhao is with the Department of Applied Biology and Chemical Technology, The Hong Kong Polytechnic University. Email: q.zhao@polyu.edu.hk.
- Le Song is with GenBio AI and Mohamed bin Zayed University of Artificial Intelligence. Email: le.song@mbzuai.ac.ae.

(Corresponding authors: Dr. Wenqi Fan and Prof. Qing Li.)

1. In this work, the terms "amino acid" and "residue" are used interchangeably.

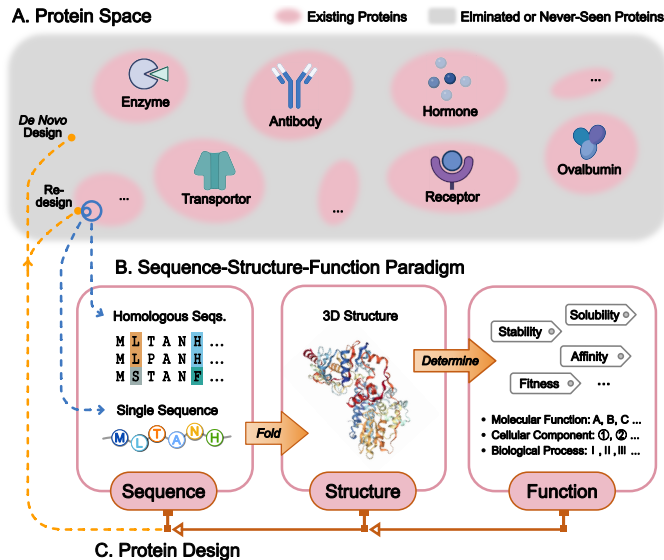


Fig. 1: Illustration of the Evolution and Sequence-Structure-Function Relationships. (A) The arrangement of amino acids forms a vast space of possible protein sequences. However, only a few proteins can survive through millions of years of evolution. (B) Valid amino acid sequences would fold into stable 3D structures and carry out proper functions. (C) Information flow within the sequence-structure-function paradigm can be leveraged in reverse, leading to the optimization of existing proteins or *de novo* protein design oriented by desired functions.

protein modeling methods excel in their own specific tasks, they still have certain limitations. To be specific, traditional AI4Protein models, even some of their protein representation learning methods, cannot sufficiently understand the critical "grammar" and deep "semantics" within protein language. This limitation arises from their inferior capabilities in sequence processing and the integration of world knowledge. Meanwhile, most existing protein modeling methods are designed for specific tasks and lack the generalization capability needed for multiple and even unseen tasks in the training stage.

Recently, Large Language Models (LLMs) have represented the forefront and peak of Artificial Intelligence, characterized by millions to billions of parameters and large-scale training based on extensive datasets. As is widely recognized, LLMs like BERT [13], T5 [14], GPT series [15, 16, 17] and LLaMA [18] have achieved remarkable success and boosted the comprehensive advancements in various research fields, such as Natural Language Processing (NLP) [19], Recommender Systems [20, 21], Healthcare [22], and more. Essentially, LLMs spark with two unprecedented superiorities. *First*, through large-scale pre-training from the extensive open-world knowledge, LLMs acquire foundational emergent capabilities, particularly in language understanding and generation. For text, a typical kind of sequence data, LLMs excellently grasp the arrangement patterns of tokens that conform to the basic grammar and express high-level semantics. *Second*, LLMs exhibit unprecedented generalization capability with the support of techniques like probing, fine-tuning, prompting, etc. They are able to address a wide array of problems, moving beyond merely improving the prediction performance on individual specific

tasks. Thus, building on the successful experience in natural language processing, LLMs have shown great potential for deciphering other "language sequences." Specialized LLMs have been developed to process genomic sequences [23], chemical molecules [24], and especially — the proteins.

By combining powerful LLM techniques with abundant protein data, various **protein Language Models (pLMs)** are proposed to sufficiently grasp the foundational protein knowledge. More specifically, sequence-only pLMs (e.g., ESM-2 [25], ProtGPT2 [26], xTrimoPGLM [27], ESM-MSA-1b [28]) capture the valid AA arrangement patterns that have emerged over the course of evolution. pLMs that incorporate explicit structure and function information (e.g., SaProt [29], ESM-3 [30]) are effectively enhanced in protein understanding and generation. pLMs that learn natural language as well (e.g., ProLLaMA [31], BioT5 [32]) can understand a broad biomedical background and have the ability to follow textual instructions. Moreover, pLMs have been excellently generalized to the protein structure prediction, protein function prediction, and protein design problems — those studies of sequence-structure-function reasoning have ushered inspiring progress simultaneously. For example, the representations learned by pLMs have empowered the accurate and fast inference from a single protein sequence to a 3D structure at atomic resolution [25]; the "pLM encoder - LLM decoder" framework unifies multi-task protein function prediction in a question-answering form [33]; the autoencoding and autoregressive pLMs well assist in protein redesign [34] and *de novo* protein design [35], respectively. Taking all of those booming research trends into account, it's imperative to conduct a thorough review of pLMs and the latest advances in computational protein science promoted by them.

In this survey, we provide a systematic overview of computational protein science empowered by LLM techniques, so as to help researchers with an AI or biology background quickly understand relevant developments and gain insights. In particular, the rest of this paper is organized as follows: Section 2 presents a background of protein data profiles, AI for protein, and large language models. Section 3 categorizes the existing pLMs into sequence-based ones, structure-&-function enhanced ones, and multimodal ones. Section 4 summarizes the utilization and adaptations of pLMs by considering the pending problems of protein structure prediction, protein function prediction, and protein design. Section 5 introduces some biomedical applications, including antibody design, enzyme design, and drug discovery. Section 6 provides a discussion on the current challenges and potential future directions. Our main contributions can be summarized as follows:

- We conduct a comprehensive literature review on computational protein science in the era of LLMs, covering the foundational pLMs, the utilization and adaptation of pLMs, and some biomedical applications.
- We present a systematic categorization for pLMs, emphasizing their learned knowledge of protein sequence, protein structure & function, and external languages, and outlining the mainly employed approaches in building pLMs.
- We delve into the utilization and adaptation of pLMs, highlighting their impacts on the protein sequence-structure-

function reasoning problems, and introducing the typical technical strategies.

- We describe some of the latest applications of pLMs, and discuss the prospective future directions in this field.

To date, there are a limited number of existing surveys related to this topic. Regarding foundation models, Zhang et al. [36] enumerate existing pLMs and categorize them by network architectures, and Hu et al. [37] introduces the connection and progression between pLMs and structure prediction. Focus more on specific protein modeling problems, Unsal et al. [38] implement a benchmark for functional property learning of proteins based on language models, while Ferruz et al. [1] and Rufflo et al. [39] both report the controllable protein design with language models. In summary, they have not comprehensively encompassed the contents of foundation pLMs, downstream generalizations, and real-world applications, and their categorizations can be further improved by balancing the challenges in scientific explorations and technical strategies. Therefore, a comprehensive survey with systematic categorizations in this rapidly developing field is still in demand.

## 2 BACKGROUND

In this section, we provide a brief review of the relevant background, including the biological basis and data profiles of proteins, the AI for protein studies, and large language models (LLMs).

### 2.1 Biological Basis and Data Profiles

Understanding the synthesis, evolution, structure, and function of proteins would pave the way for advances in biology and medicine. Over the years, increasing attention has been paid to studying proteins for various tasks in science. For example, many scientists from the wet lab have conducted comprehensive wet experiments through advanced biotechnology (e.g., mass spectrometry [40], X-ray crystallography [41], electron microscopy [42], and deep mutational scanning [43]), significantly contributing to the accumulation of a substantial volume of high-quality data [44, 45, 46]. Thus, we start by providing some background biological knowledge and representative data formats in Figure 2.

In the wild, protein synthesis involves two processes: transcription and translation. Genetic information is transcribed from protein-coding genes into messenger RNAs (mRNA), which are subsequently translated into proteins [47]. Specifically, the mRNA sequence is composed of four types of nucleotides, i.e., adenine (A), uracil (U), cytosine (C), and guanine (G), while a codon refers to a subsequence of three consecutive nucleotides, such as "AUG", "AAG", etc. According to the genetic code [48], the 20 standard AAs are encoded by 61 codons. For example, Methionine (abbr. as Met or M) is specified by the codon "AUG", while Lysine (abbr. as Lys or K) can be translated from the codons "AAA" or "AAG". Through the sequential translation of codons and post-translational modifications, a protein is primarily constituted by a sequence of amino acids, which is also called the polypeptide chain. We can denote it as  $X = [x_1, x_2, \dots, x_L]$ , where  $L$  refers the sequence length, and  $x_i$  ( $1 \leq i \leq L$ ) represents each individual amino acid.

The total number of possible amino-acid sequences can easily exceed orders of billions, yet only a minute fraction has existed on earth since the origin of life [5]. Sequences that fail to fold into stable structures and exhibit effective functionality are more likely to be eliminated in the course of evolution, while the survived natural protein sequences would reflect evolutionary favored patterns. To be specific, evolution gives rise to the emergence of protein families that share significant structural and functional similarities [49]: Random genetic mutations occur in an ancestral protein over time, and certain descendant variants survive through the process of natural selection; Beneficial mutations accumulate in the ancestral protein within the population across generations, which leads to the adaptation and divergence of the protein family. In computational analysis, it's a common practice to retrieve homologous sequences and create multiple sequence alignment (MSA) to incorporate the prior knowledge of protein evolution. Concretely, an MSA comprising  $M$  sequences of length  $L$  can be encoded as a matrix  $X \in \mathbb{R}^{M \times L}$ , where each entry  $x_{m,l}$  ( $1 \leq m \leq M, 1 \leq l \leq L$ ) represents the amino acid identity of sequence  $m$  at position  $l$ . By comparisons, we can identify the conserved positions that remain unchanged, the variable positions that allow different kinds of mutations, and the co-evoluted positions that are changed synchronously as constrained by structural contacts.

After the synthesis of polypeptide chains, protein folding is a fundamental physical process that transforms the unstable sequential conformations into more ordered three-dimensional structures [50]. The primary structure of a protein refers to the linear arrangement of AAs in the polypeptide chain, which holds great importance in encoding the protein's 3D structure and biological function from the very beginning. Then, the secondary structure refers to highly regular local sub-structures along the polypeptide chain, typically including  $\alpha$ -helices and  $\beta$ -sheets. The tertiary structure refers to the 3D structure formed based on a single polypeptide chain, where a series of  $\alpha$ -helices and  $\beta$ -sheets are folded compactly through non-specific hydrophobic interactions. The quaternary structure describes the 3D structure formed by aggregating at least two individual polypeptide chains, which function together as multimers. Protein Data Bank (PDB) [45] is an authoritative database that collects experimentally determined protein structures, while "*pdb*" and "*mmcif*" are widely acknowledged file formats that provide detailed descriptions of the 3D structures of proteins. A structure file comprehensively records the  $(x, y, z)$  coordinates for each atom constituting all amino acids. Then, the structure of a protein can be redescribed by executing various data processing schemes. For example, we can calculate the relative distance, angle, and direction among atoms or AAs, and build a graph that represents AAs as nodes and their spatial relations as edges [51, 52].

A fundamental principle in molecular biology asserts that the structure of a molecule inherently determines its function [4]. For example, antibodies are the Y-shaped proteins employed by the immune system to specifically identify and neutralize pathogenic organisms, like bacteria and viruses. The two arms in the "Y"-like structure are crucial as hosting antigen-binding sites that recognize and bind to antigens. Then, the detailed conformation of antigen-binding sites further determines the antibody's functional affinity

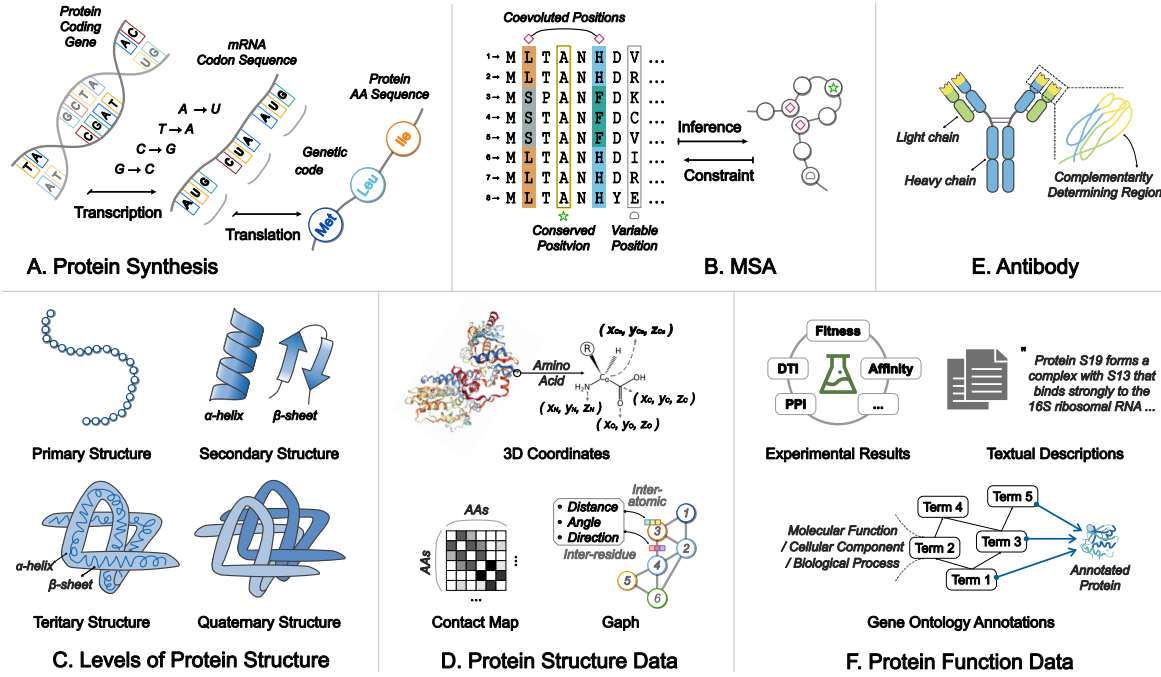


Fig. 2: Biological basis and data profiles. (A) Protein synthesis mainly involves the transcription of protein-coding genes to mRNAs and the translation of codon sequences to AA sequences. (B) Multiple Sequence Alignment (MSA) contains the evolutionary prior knowledge of proteins. Conserved positions are interpreted as core AAs for protein structure, as no changes have been allowed throughout the evolutionary process. Pairs of coevolved positions indicate the spatial contacts of AAs, since mutations would occur and act synergistically to preserve the structure stable and unchanged. (C) Protein structure exhibits hierarchical organization. (D) Protein structure can be described in several forms. 3D coordinates of atoms trustfully record the experimentally determined protein conformation. A 2D distance map conveys the proximity between all possible AA pairs. Furthermore, we can build specific graphs to describe detailed structural characteristics, where the interatomic or inter-residue distances, angles, and directions are encoded as node and edge features. (E) An antibody is a Y-shaped protein composed of two heavy and two light chains. At the top of the "Y"'s arms, complementarity-determining regions (CDRs) are polypeptide segments that make up the antigen binding site. (F) Protein function is described in multiple formats, such as lab-generated labels, Gene Ontology annotations, and textual documents.

for clearing specific pathogens. Beyond the antibodies, there are broader facts that demonstrate the complexity of protein function mechanisms. Structural differences among proteins correlate strongly with functional diversity, and even homologous proteins with similar structures can exhibit significant variations in the fitness landscape [53]. Considering the great complexity, protein function is presented by a wide range of data formats, and we summarize them as experimental results, manual annotations, and textual descriptions:

- First, scientists conduct laboratory experiments to measure functional properties (e.g., fitness, stability), identify functional sites (e.g., signal peptides, B-cell epitope), and detect biomolecular interactions (e.g., drug-target interaction, protein-protein interaction). All kinds of results are collected as labels in datasets.
- Second, proteins can be annotated with authoritative classes. Gene Ontology (GO) knowledgebase [46] organizes a hierarchical framework of terms that describe various classes of molecular functions, cellular components, and biological processes. Each protein is manually tagged with multiple GO terms, facilitating a granular and systematic exploration of protein functions. Besides, Enzyme Commission (EC) number is an authoritative taxonomy that assigns a unique code to each enzyme based on their catalyzed chemical reactions. Like "EC 2.7.6.1", an enzyme code is composed of a beginning string "EC" and four numbers separated by dots, denoting a progressively finer

classification of the enzyme.

- Third, academic publications serve as an essential resource, providing comprehensive empirical evidence and theoretical analyses for understanding protein functions.

## 2.2 AI for Protein Science

There are essential protein modeling problems centered in the sequence-structure-function paradigm, including *protein representation learning*, *structure prediction*, *function prediction*, and *protein design*. In recent years, artificial intelligence has significantly advanced these studies.

As one of the most representative techniques in protein modeling, protein representation learning aims to extract latent knowledge of proteins from extensive data and encode individual proteins into vector representations. Then, the learned protein knowledge can be used for various downstream applications. As summarized by Wu et al. [54], existing protein representation learning methods can be categorized into three groups: sequence-based protein encoders, structure-based protein encoders, and sequence-structure co-modeling methods. First, UniRep [7], CARP [55], and ProtTrans [56] train LSTM, CNN, and Transformer models on single AA sequences, while ESM-MSA-1b [28] propose customized axial Transformers for multiple sequence alignments (MSAs). These encoders are mainly pre-trained using the mask language modeling objective.

Second, GearNet [57] and GVP-GNN [51] utilize message-passing graph neural networks (GNNs) or geometric-aware GNNs on protein structure graphs. These encoders learn by contrasting sampled structures and incorporating specific supervision on structural characteristics. Third, LM-GVP [58] and ESM-GearNet [59] introduce novel strategies that combine sequence models and GNNs by modifying the network architectures and learning objectives. It is worth noting that protein language models (pLMs) are closely connected with protein representation learning. Certain pLMs, particularly encoder-only models, can be acknowledged as protein representation methods. In addition, pLMs are commonly integrated within protein representation methods as a crucial component, working alongside modules that extract structural or functional features from proteins.

Protein structure prediction is designed to infer the 3D structure of a protein according to its amino acid sequence. At a lower resolution, protein structure prediction is defined as secondary structure prediction or contact map prediction: An AA sequence is mapped to a series of secondary structure elements (i.e.,  $\alpha$ -helix,  $\beta$ -sheet, or other conformation); Each pair of residues in the sequence is mapped to whether they are in contact. Furthermore, AI models have made substantial progress in predicting protein structures at the atomic resolution, i.e., inferring accurate 3D coordinates. Prominent methods, such as AlphaFold2 [8], RoseTTAFold [9], ColabFold [60], etc., leverage evolutionary information to achieve unprecedented performance. These methods typically involve a two-step process: first, they search for homologous multiple sequence alignments (MSAs) and structure templates according to the input sequence; then, they generate coordinates for all atoms through well-designed networks that collectively analyze relationships within and between the 1D sequence, 2D contact, and 3D conformation. Despite achieving near experimental accuracy, these methods are reliant on MSAs and may use known structures from homologous proteins as references, which deviates slightly from the goal of capturing the underlying rules that govern polypeptide folding. In order to overcome this deficiency, efforts have been dedicated to single-sequence protein structure prediction. Certain pLM-based methods, such as ESMFold [25], trRosettaX-Single [61], and HelixFold-Single [62], demonstrate competitive performance comparable to AlphaFold2 and RoseTTAFold while free of the prior co-evolution information.

Protein function prediction aims to conclude the functional knowledge of proteins based on their AA sequences or 3D structures. Given the diverse nature of protein functions, protein function prediction encompasses various concrete tasks. First, when considering each input protein as a whole, AI techniques are employed to predict experimentally measured properties (e.g., stability, fluorescence, fitness) [63] or manually curated annotations (e.g., gene ontology annotations, enzyme classification numbers) [10, 64]. Second, proteins perform specific functions through the finer functional sites, such as signal peptides, complementarity-determining regions, ligand-binding domains, etc. Consequently, AI models are developed to identify whether each residue belongs to a specified functional region [65, 66]. Third, proteins generally don't function in isolation but rather interact with biomolecules (e.g., DNA, RNA, drug, another protein)

to carry out their functions. AI models are designed to determine whether interactions occur or to reveal interaction details based on pairwise protein-molecule inputs [67, 68]. In addition, a cutting-edge research spot involves enabling multi-task learning within a unified framework, which eliminates the need for numerous individual models. Efforts are made to develop ChatGPT-like systems that allow users to upload proteins and engage in question-answer conversations to gain insights [33, 69, 70].

Protein design aims to create new proteins with functions surpassing existing proteins by discovering optimized or novel sequences [71]. It is required to precisely identify specific regions within the vast protein sequence space that give rise to desired functions. Depending on whether starting from known proteins or scratch, relevant studies can be classified into two groups: protein redesign and *de novo* protein design. During protein redesign, mutations are introduced to existing proteins to enhance their functional properties. Therefore, limited protein space surrounding existing proteins is explored. Autoencoding pLMs trained by masked language modeling excel at predicting highly probable mutations favorable in evolution, inspiring "new family members" that are directionally evolved [72, 34]. Dissimilarly, *de novo* protein design involves autonomously discovering optimal amino acid sequences that satisfy a given design objective, enabling the exploration of regions not previously seen in evolutionary history within the extensive protein space. Most *de novo* protein design methods divide the workflow into two stages [73]: 1) Structure generation models (e.g., RFDiffusion [74] and Chroma [75]) create the protein's structural backbone without a defined sequence; 2) Inverse folding models (e.g., ESM-IF [76] and ProteinMPNN [77]) recover effective sequences based on the backbone atom coordinates. Besides, certain studies consider *de novo* protein design as a single-step computational task, i.e., protein sequence generation. Autoregressive pLMs (e.g., ProGen [35] and ProtGPT2 [26]) and diffusion frameworks (e.g., EvoDiff [78]) can generate protein sequences that exhibit similar physicochemical properties to natural proteins but have low homology to existing proteins.

## 2.3 Large Language Models (LLMs)

Recently, the development of LLMs has unleashed a great surge in AI [17, 19]. Typically, LLMs with million or billion-level parameters are pre-trained on extensive data, demonstrating surprising abilities in various NLP tasks like language understanding, text generation, etc. Due to their powerful understanding and generation capability, LLMs have been widely adopted in various fields, including data mining [20, 79], healthcare [22] and molecule science [24, 80], by fine-tuning or prompting with domain-specific datasets. According to the architecture and functionality, LLMs are categorized into three primary types: *Autoencoding models*, *Autoregressive models* and *Sequence-to-Sequence models*.

As a typically Autoencoding model, BERT [13] transforms input text into a high-dimensional latent space to capture the contextual semantic information of the text. One of the key features of Autoencoding models is their bi-directionality nature, which allows them to process both the preceding and following context. This feature enables Autoencoding models



to better understand and represent the input data, which is important in tasks like sentiment classification and language translation. Different from Autoencoding, Autoregressive models such as GPT [15] family processes the input text in a left-to-right manner. They generate the next token by predicting based on the context of previous tokens. This autoregressive generation fashion allows them to generate coherent and contextually relevant text in tasks such as creative writing and code generation. In addition, Sequence-to-Sequence models like T5 [14] take conditional generation as the core idea. Specifically, they use an encoder to capture the meaning of the input text and a decoder to generate the output based on encoding information. This architecture makes Sequence-to-Sequence models a unified framework to handle various NLP tasks flexibly. It is worth noting that the performance of LLMs follows a scaling law [81]. Specifically, the performance of LLMs has been observed to scale predictably with increases in model size and the amount of training data [17].

One typical method for adopting LLMs to various downstream predictions is fine-tuning the LLMs on specific datasets. That is to say, the pre-trained parameters are subsequently trained to learn more specific tasks. However, full model fine-tuning still requires considerable data samples, costs extensive computing resources, and takes quite some time. To overcome these limitations, recent studies explore Parameter-Efficient Fine-Tuning (PEFT) to adopt partial parameters instead of fine-tuning the LLMs extensively [82, 83]. For example, Low-Rank Adaptation (LoRA) [82] allows for the effective adaptation of LLM by only updating a few low-rank matrices, thus substantially decreasing the number of parameters while preserving high-performance levels.

To further reduce the reliance on training data, prompt learning recently emerged as a popular paradigm for adopting LLMs to various tasks, with the input carefully designed to guide LLMs in generating the desired output without intensive parameter updates. In-context learning is a commonly used prompt learning method that provides a few task demonstrations within the prompt to instruct the LLMs to perform downstream tasks. Furthermore, Chain-of-Thought (CoT) [84] is another specific prompting technique. The key idea of CoT prompting is to annotate intermediate reasoning steps into the prompt to enhance the reasoning ability of LLMs. While these prompting methods achieve great success, such manually designed prompts typically face discrete optimization challenges, such as laborious trial and error in finding suitable prompts. To solve this challenge, soft prompt tuning [85] is introduced, where prompts are not fixed text but rather continuous, trainable embeddings, allowing a more nuanced and flexible prompt design.

Furthermore, recent studies have expanded the utility of LLMs by integrating them with other modalities or external knowledge bases to enhance their performance and versatility. For example, models like CLIP [86] utilize contrastive learning to comprehend images through textual descriptions, showcasing exceptional abilities in cross-modal understanding. PaLM-E architecture [87] integrates the specific encoder and the uni-decoder hierarchically with the help of an internal projector [88], enabling it to generate contents based on the understanding of visual and language inputs. Similarly, BLIP-2 [89] introduces a frozen image

encoder with a pre-trained text decoder, where a lightweight Querying Transformer (Q-Former) module is developed to bridge the modality gap. Furthermore, without updating the LLM backbone, Retrieval-Augmented Generation (RAG) techniques have emerged as an effective approach to enhance the understanding and generation capabilities of LLMs by retrieving external knowledge database [90].

### 3 PRE-TRAINED PROTEIN LANGUAGE MODELS

In recent years, LLM techniques have been employed not only in processing natural languages but also in deciphering various domain-specific "languages". Substantial advancements have been made in protein language models (**pLMs**). In this section, we systematically review existing pLMs, categorizing them as sequence-based, structure-and-function-enhanced, and multimodal models.

#### 3.1 Sequence-based pLMs

General LLMs capture the interdependencies among sub-word tokens and acquire a profound understanding of the grammar and semantics of text. Similarly, sequence-based pLMs capture the mutual dependencies among amino acid (AA) tokens, extract the favorable sequence patterns, and grasp implicit structural and functional information. Sequence-based pLMs can be further distinguished as single-sequence-based ones and multiple-sequence-based ones. The former describes each protein by the corresponding AA sequence, while the latter possesses an idea of retrieval augmentation, describing each protein with multiple related sequences in evolution or synthesis. Table 1 presents a comprehensive summary, outlining the input data, network architecture, and pre-training objective of each pLM.

##### 3.1.1 Single-Sequence-based pLMs

By taking each amino acid sequence as a sentence, the valuable practice of building general LLMs has been extended to the development of pLMs. We have witnessed the emergence of autoencoding pLMs in the BERT-style, autoregressive pLMs following the GPT approach, and sequence-to-sequence pLMs resembling T5 or GLM.

As shown in Figure 3-1, **Autoencoding pLMs** employ the Transformer encoder with bidirectional attention and undergo pre-training via the masked language modeling (MLM) objective. They are skilled in encoding the context of protein sequences into informative representations. ESM-1b [91] is a pioneering study that executes self-supervised learning on 250 million protein sequences spanning evolutionary diversity. Technically, ESM-1b employs the architecture and pre-training procedure of RoBERTa [92] almost as-is, thereby driving the model to extract the residue-residue dependencies latent in the extensive sequence data. Notably, it is revealed that the structural residue-residue contacts can be inferred from the pLM-produced protein sequence representations by linear projections. The significant correlation between the actual distance map and the extracted inter-residue dependencies underscores a valuable consistency between the nature of proteins and the framework of language models. With a similar idea, ProtTrans [56] involves several other representative autoencoding LMs to learn the protein

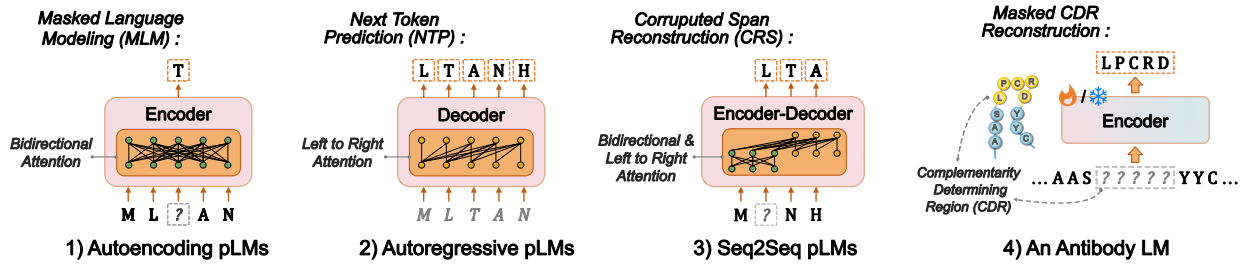


Fig. 3: Typical **single-sequence-based pLMs**. 1-3) When considering individual amino acid sequences as "sentences", pLMs follow the general approaches of autoencoding, autoregressive, and sequence-to-sequence as well. 4) Masked CDR reconstruction is a novel pre-training objective that incorporates the inherent characteristics of antibodies into mask language modeling.

language, resulting in ProtBERT, ProtAlberty, and ProtElectra. In transfer learning evaluations, those single-sequence-based presentations perform competitively with classic methods that take in multiple sequence alignments. This finding implies that pLMs have captured some of the grammar of the protein language, even without direct exposure to evolutionary information.

Subsequently, to explore the boundaries of sequence-based protein representation learning for optimal performance, the next-generation study of ESM-1b, i.e., ESM-2 [25], implements the scaling law of LLMs. With the scale of parameters increasing from 8 million to 15 billion, significant improvements are observed in the fidelity of protein sequence modeling, leading to the emergence of protein structure knowledge at the atomic level within the learned representations. As of December 2024, the ESM team has made the latest generation, ESM Cambrian (ESM C) [93], available. It is claimed that ESM C establishes a new state-of-the-art performance for protein sequence modeling, achieving linear scaling across orders of magnitude in parameters. With improved training efficiency, each ESM C model could match or even greatly exceed the performance of previous larger ESM-2 models.

Besides, efforts have been made to improve autoencoding pLMs from different points of view. To better capture co-evolutionary information in inter-residue co-variation, He et al. [94] introduce a novel pre-training objective called pairwise masked language modeling (PMLM). For extrapolation of longer proteins and protein complexes, LC-PLM [95] is developed using a novel compute-efficient network architecture called bidirectional Mamba with shared projection layers (BiMamba-S) [96, 97]. To mitigate computational resource requirements, ProtFlash [98] employs a mixed chunk attention strategy with linear complexity, while DistillBERT [99] is proposed as a distilled variant of ProtBERT. To find a balance between performance and efficiency in the trend of scaling up pre-training, AIDO.Protein [100] explores the Mixture-of-Experts (MoE) model in computational protein science for the first time. It serves a critical link in the AI-Driven Digital Organism (AIDO) [101] system, which aims to integrate multi-scale foundation models to predict, simulate, and program biology at all levels. Besides, DPLM [102] is a versatile protein language model under a discrete diffusion framework [103]. While maintaining the bi-directional receptive field, the key ingredient of autoencoding modeling, an additional diffusion pre-training process makes DPLM possess a strong and scalable generative power.

As shown in Figure 3-2, **Autoregressive pLMs** leverage Transformer decoder to process the protein sequence left-to-right and conduct pre-training using the next token prediction (NTP) objective. They are good at generating protein sequences that adhere to the patterns favored by evolution. Building upon the approach of GPT-2, ProtGPT2 [26] is capable of generating protein sequences that exhibit natural amino acid propensities yet relate distantly to naturally occurring sequences. Across the progression from GPT-2 to GPT-3 and subsequent iterations, general LLMs have consistently exhibited a trend of scaling up, which also applies to pLMs. RITA [104] encompasses a suite of models ranging from 85 million to 1.2 billion parameters, which are trained on over 280 million protein sequences. ProGen2 [105] models are enlarged to 6.4 billion parameters and are trained on an extensive collection of over one billion protein sequences. Moreover, it is observed that even the largest ProGen2 model still exhibits underfitting, indicating the potential for further improvements in capturing the intrinsic distribution of natural protein sequences.

As illustrated in Figure 3-3, **Sequence-to-sequence pLMs** employ either the encoder-decoder or non-causal decoder architecture and undergo pre-training with the corrupted spans reconstruction (CSR) objective. These models combine autoencoding and autoregressive traits, can encode the input sequence, and perform conditional generation accordingly. In ProtTrans [56], ProtT5 is a collection of pLMs inspired by the original T5 series. While scaling up to 3 billion and 11 billion parameters, ProtT5 models produce highly informative protein representations that outperform the smaller autoencoding pLMs such as ProtBERT. Following ProtTrans, Ankh [106] is another empirical study that explores data-efficient, cost-effective, and knowledge-guided optimization of pLMs. Using ProtT5 as the baseline, Ankh comprehensively investigates potential factors that could affect the performance of pLMs, covering more than twenty ablation experiments that compare detailed strategies in token span corruption, encoder-decoder architecture, and data sources. Besides, xTrimopGLM [27] is also proposed to process the protein representation and protein generation uniformly. xTrimopGLM leverages General Language Model (GLM) [107] as its backbone architecture and explores the joint optimization of MLM and CSR objectives, resulting in successful pLM pre-training at the scale of approximately 100 billion parameters and 1 trillion tokens as never before.

In contrast to general proteins, antibodies (i.e., immunoglobulins) are specialized Y-shaped proteins with

immune functions. As illustrated in Figure 2-E, antibodies are characterized by the two heavy chains and two light chains in their sequences. Then, complementarity-determining regions (CDRs) are vitally important components within antibody sequences that dictate the specific antigen-binding sites. Considering the critical medical significance and distinctive sequence organization of antibodies, **antibody language models** have been developed specifically. In early attempts, AntiBERTy [108], AntiBERTa [109], and AbLang [110] employed the classic autoencoding approach of BERT or RoBERTa to decipher antibody sequences. Furthermore, the pre-training objective of antibody LMs can be enhanced by incorporating relevant biological knowledge. As shown in Figure 3-4, AbBERT [111] models are pre-trained by reconstructing masked CDR spans, thereby acquiring antibody-specific insights. ReprogBert [112] investigates cross-language adaptation using limited data, reprogramming an English BERT model into an antibody LM by training amino acid vocabulary embeddings with the masked CDR reconstruction objective. Besides, autoregressive and sequence-to-sequence antibody LMs have also been developed. For example, IgLM [113] enables the generation of antibody sequences conditioned on prefix tag tokens that indicate the chain type and the origin species, paired-IgGen (p-IgGen) [114] is a generative LM pre-trained on both paired and unpaired heavy-light chain antibody sequences, and pAbT5 [115] exhibits the "translation" capability that generates light-to-heavy-chain or heavy-to-light-chain sequences. As learning the distinct language of antibodies, these models effectively contribute to the controllable design of antibodies.

### 3.1.2 Multiple-Sequences-based pLMs

Retrieval is a fundamental data mining technique designed to understand input queries and extract relevant information to assist in analysis [90]. In computational protein science, retrieval has been a well-established concept for decades. Many significant analyses are conducted based on multiple biologically related sequences, aiming to leverage prior evolutionary knowledge, such as the co-variation between residues. Driven by the concept of retrieval and the practice in large-scale pre-training, multiple-sequences-based pLMs have been developed as well.

Traditional searching tools like HHblits [116] and MM-seqs2 [117] are widely employed to retrieve and align protein families from extensive databases. The resulting multiple sequence alignments (MSAs) can greatly contribute to investigating the evolutionary relationships between proteins. As shown in Figure 2-B, MSAs facilitate the identification of conserved, variable, or co-evolution regions within the protein sequences, offering valuable insights into their structural and functional implications. ESM-MSA-1b [28] is the first pLM designed to operate on MSAs. As illustrated in Figure 4-1, ESM-MSA-1b incorporates bidirectional tied-row attention and column attention within each MSA Transformer block, thereby capturing the interdependencies among residues and across sequences. In evaluations, ESM-MSA-1b demonstrates great performance in unsupervised protein structure learning, surpassing single-sequence-based pLMs while with far fewer parameters. These findings meet the motivation that structural constraints can be effectively inferred from the patterns of protein sequences.

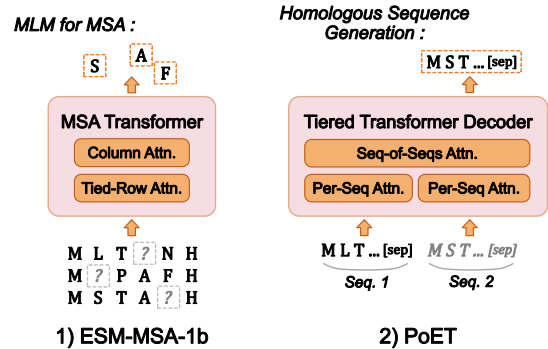


Fig. 4: Typical **multiple-sequences-based pLMs**. 1) ESM-MSA-1b, a representative MSA-based pLM, incorporates bidirectional tied-row attention and column attention within each MSA Transformer block, thereby capturing co-evolution features within the 2D input. 2) PoET is an autoregressive model specifically designed to learn the distribution over protein families. It accepts multiple sequences as input without the need for alignment and can generate sets of homologous proteins.

Moreover, as well-established protein structure prediction methods (e.g., AlphaFold2) rely heavily on MSAs and exhibit inadequacy when confronted with orphan proteins lacking adequate homologous sequences, there is an increasing demand to generate "pseudo" proteins based on existing MSAs for augmentation. MSA2Prot [118] is an MSA-to-protein Transformer that can generate individual protein sequence variants based on a learned encoding of input MSA. Then MSA-Augmenter [119] leverages the tied-row & column attention mechanism to generate an arbitrary number of co-evolutionary sequences. Besides, MSAGPT [120] involves a flexible MSA decoding framework. By employing a 2D positional encoding scheme that describes the complex evolutionary patterns, MSAGPT empowers the general 1D Transformer decoder for MSA processing.

However, from another point of view, the efficacy of MSA-based pLMs is usually hindered by the inadequacy of fundamental alignment algorithms, which can introduce errors such as long insertions and gappy regions. Even worse, the alignment errors tend to increase, and the computational efficiency decreases when handling longer sequences. To avoid these problems, certain pLMs are designed to operate on multiple sequences without the need for alignment. As shown in Figure 4-2, PoET [121] contains a tired sequence-of-sequences causal attention mechanism to generate sets of homogeneous protein sequences. Besides, inspired by the recent advancement of Mamba [96], ProtMamba [122] efficiently handles significantly long contexts, even encompassing hundreds of sequentially concatenated protein sequences.

In addition to retrieving homogeneous AA sequences, proteins can be supplemented by other biological sequences involved in the synthesis process. As introduced in Figure 2-1, each AA sequence naturally originates from an mRNA codon sequence, where genetic information can be exploited to build enhanced pLMs. CaLM [123] is designed to generate representations of codon sequences, providing valuable signals for protein engineering applications. cdsBERT [124] introduces a pipeline to enhance the capability of pre-trained ProtBERT models by introducing codon awareness. Furthermore, post-translational modifications (PTMs) represent a



covalent process that modifies proteins after their synthesis, playing a vital role in increasing proteins' structural and functional diversity. PTM-Mamba [125] stands as the first PTM-aware pLM, encompassing representations of both amino acids and PTM tokens.

### 3.2 Structure and Function Enhanced pLMs

As sequence-based pLMs demonstrate the ability to capture implicit structural and functional semantics from protein sequences through large-scale pre-training, the further integration of explicit knowledge can enhance their understanding of proteins at a more comprehensive level. In this subsection, we present the recent advancements in constructing structure-and-function-enhanced pLMs. We explain the data form of protein structure and function individually, and introduce the corresponding incorporation methodologies. Relevant contents are also summarized in Table 2.

#### 3.2.1 Structure Enhanced pLMs

As of December 2024, the Protein Data Bank (PDB) [45] has accumulated a vast repository of over 220 thousand experimentally determined protein structures. Furthermore, AlphaFold Protein Structure Database offers more than 200 million reliable structure prediction results. These invaluable resources establish explicit mapping relationships between sets of protein sequences and their 3D structures, serving as a robust data foundation for structure-enhanced pLMs. Figure 5 presents three typical cases of integration of structural information through pre-processed structural features, the structural graph, or discrete tokens.

Initially, each PDB file comprehensively documents the atom-scale 3D coordinates of a protein structure, and we can extract a wide variety of structural characteristics for more targeted learning. For example, PeTriBERT [142] incorporates additional position encodings for residues, including their central coordinates and rotation parameters. ESM-S [143] enhances ESM-2 by integrating a new training task, remote homology detection, which seeks to identify proteins with similar structures but low sequence similarity. METL [144] is trained to infer protein biophysical attributes, which are calculated from static structures through Rosetta [145].

Moreover, protein structures can be represented as graphs, where the nodes refer to amino acids, and the edges connect AAs that are close in distance. There is an important conjunction between the nature of graph and the principle of language models: graph represents relationships between nodes, while language models learn relationships between tokens. It is revealed that residue-residue contacts can be predicted from the attention maps of sequence-based pLMs [146]. Inversely, the structural graph can be incorporated into the attention map of pLMs in the construction of structure-enhanced pLMs. Additional GNN encoders (e.g., GVP [51], GearNet [57]) that represent structural inter-residue relations are integrated together with the attention mechanism of pLMs. For instance, PST [147] refines ESM-2 by attaching a novel structure extractor into each attention module, LM-Design [148] introduces a lightweight structural adapter into ESM-1b after the final Transformer layer, and ProseLM [149] introduce well-designed adapters that update the outputs of the attention & feed-forward operations of each ProGen2

layer in fusing encoded structural features and low-rank language model embeddings.

In recent years, vector quantized-variational autoencoder (VQ-VAE) [150, 21] has emerged as a flourishing technique that encodes the 3D protein structure into discrete tokens, representing the local geometric conformation of each amino acid. Technically speaking, VQ-VAE employs an encoder and a decoder to map data between the real and hidden space, and learns a codebook to quantize continuous latent representations into discrete ones. As the VQ-VAE tokenizer compresses intricate continuous data into a limited number of discrete latent representations, there is a basic consensus that a small codebook size usually leads to more information loss. In retrospect, Foldseek [151] pioneeringly leverages VQ-VAE to establish the 3D interaction (3Di) alphabet for protein structure. The maximally evolutionary conserved structure states are represented by twenty 3Di tokens, which facilitate fast and accurate protein structure searching. While Foldseek's 3Di tokens have demonstrated robust performance in protein retrieval, the coarse-grained nature of the small codebook still limits its structure reconstruction ability. To make up for this weakness, more protein structure tokenizers have been proposed, featuring inventive designs of the encoder, quantization, and decoder modules, along with the critical exploration of codebook size. ProTokens [152], FoldToken series [153], and AIDO.st [154] are all promising approaches to distill protein structure tokens, aiming to achieve the balance among varied factors, including the compression efficiency, retrieval ability, reconstruction ability, downstream usability, etc.

With protein structure tokenization methods, each 3D structure can be described as an array of protein structure tokens, allowing for seamless integration with language models. For instance, ProstT5 [155] is developed by incrementally training ProtT5 to translate between the structure 3Di tokens and the sequence AA tokens. SaProt [29, 156] leverages a Structure-Aware (SA) vocabulary that integrates AA tokens and 3Di tokens to effectively represent proteins in both perspectives of primary and tertiary structures. ProSST [157] involves a well-designed sequence-structure disentangled attention to learn combined the relationship between protein AA sequences and protein structural token sequences. In protein representation learning benchmarks, those models excel in various downstream tasks [158, 159], especially in the zero-shot mutation effect prediction [160]. Besides, DPLM-2 [161] learns the joint protein sequence-structure distribution by applying discrete diffusion to the aligned AA and structure tokens. Therefore, DPLM-2 not only performs well in protein representation learning for predictive tasks but also skills in various generative scenarios, such as unconditional sequence-structure co-generation, structure-conditioned sequence generation, and so on.

Notably, ESM-3 [30] is another recent updation of the ESM series, where discrete protein semantic tokens are introduced as an essential concept. Unlike ESM-1b, ESM-2, and ESM C, which are solely learned from protein sequences, ESM-3 is designed to represent and reason over the sequence, structure, and function of proteins. The technical scheme is illustrated in Figure 5-4, ESM-3 is an all-to-all masked language model that synchronously conditions on and generates multiple separate tracks. At the input, the protein

TABLE 1: **Sequence-based pLMs**. We present the resource link, employed corpora, input data format, backbone architecture, scale of parameters, and pre-training objective for *single-sequence-based pLMs* and *multiple-sequences-based pLMs*. All abbreviations are explained in the footnote \*, and important databases are also described in footnotes 1 to 8. In the context of "Architecture", LLM or pLM names devoid of color indicate training corresponding networks from scratch, names in **blue** represent using pre-trained models while keeping parameters frozen, and names in **pink** indicate using pre-trained models with parameters trainable.

Method	Corpora	Input	Architecture	#Parameter	Objective
ESM-1b [91]	UniRef50 [44] <sup>1</sup>	AA Seq.	RoBERTa [92]	650M	MLM
ESM-2 [25]	UniRef50	AA Seq.	RoBERTa	8M ~ 15B	MLM
ESM C [93]	UniRef, MGnify [126] & JGI [127]	AA Seq.	Encoder	300M ~ 6B	MLM
ProtBERT [56]	UniRef100 [44] <sup>1</sup> / BFD [128] <sup>2</sup>	AA Seq.	BERT [13]	420M	MLM
ProtAlbert [56]	UniRef100	AA Seq.	Albert [129]	224M	MLM
ProtElectra [56]	UniRef100	AA Seq.	Electra [130]	420M	MLM
DistilProtBert [99]	UniRef50	AA Seq.	ProtBERT	230M	MLM
PMLM [94]	UniRef50	AA Seq.	RoBERTa	87M ~ 715M	MLM & Pairwise MLM
ProtFlash [98]	UniRef50	AA Seq.	Encoder w/ Linear Attn.	174M	MLM
LC-PLM [95]	UniRef50 & UniRef90	AA Seq.	BiMamba-S	130M ~ 4B	MLM
DPLM [102]	UniRef50	AA Seq.	ESM-2 [25]	150M ~ 3B	MLM & Discrete Diffusion
AIDO.Protein [100]	UniRef90 & ColabFoldDB [60] <sup>3</sup>	AA Seq.	Encoder w/ MoE [131]	16B	MLM
ProtTXL [56]	UniRef100 / BFD	AA Seq.	Transformer-XL [132]	409M / 562M	NTP
ProtXLNet [56]	UniRef100	AA Seq.	XLNet [133]	409M	NTP
ProtGPT2 [26]	UniRef50	AA Seq.	GPT-2 [16]	738M	NTP
RITA [104]	UniRef100	AA Seq.	GPT-3 [17]	85M ~ 1.2B	NTP
ProGen2 [105]	UniRef90 [44] <sup>1</sup> & BFD30 <sup>2</sup>	AA Seq.	Decoder	151M ~ 6.4B	NTP
ProtT5 [56]	UniRef50 / BFD	AA Seq.	T5 [14]	3B, 11B	MLM
Ankh [106]	UniRef50	AA Seq.	Encoder-Decoder	755M, 1.9B	CSR
xTrimoPGLM [27]	UniRef90 & ColabFoldDB	AA Seq.	GLM [107]	100B	MLM & CSR
AntiBERTy [108]	OAS [134] <sup>4</sup>	AA Seq.	BERT	26M	MLM
AntiBERTa [109]	OAS & SABDab [135] <sup>5</sup>	AA Seq.	RoBERTa	86M	MLM
AbLang [110]	OAS	AA Seq.	RoBERTa	125M	MLM
AbBERT [111]	OAS	AA Seq.	BERT	110M	Masked CDR Recon.
ReprogBert [112]	SABDab	AA Seq.	<b>BLUE</b> w/ AA Emb.	110M, 340	Masked CDR Recon.
IgLM [113]	OAS	ID Tag   AA Seq.	GPT-2	13M	CSR
g-IgGen [114]	OAS	AA Seq.	GPT-2	17M	NTP
pAbT5 [115]	OAS	AA Seq.	<b>ProtT5</b>	3B	Light-Heavy Trans.
ESM-MSA-1b [28]	UniRef50	MSA	Encoder w/ Tied Row & Column Attn.	100M	MLM
MSA2Prot [118]	Pfam [136]	MSA	Encoder w/ Axial Attn. - Decoder w/ MSA Cross Attn.	-	NTP
MSA-Augmenter [119]	UniRef50 & UniClust30 [137] <sup>7</sup>	MSA	Encoder-Decoder w/ Tied Row & Column Attn.	260M	NTP
MSAGPT [120]	UniClust30	Multi. AA Seq. + 2D PE	Decoder	2.8B	NTP
PoET [121]	UniRef50	Multi. AA Seq.	Decoder w/ Tiered Attn.	400M	NTP
ProtMamba [122]	OpenProteinSet [138] & UniClust30	Multi. AA Seq.	Mamba [139]	107M	CSR
cdsBERT [124]	CCDS [140] & Ensembl	AA Seq. & Codon Seq.	<b>ProtBERT</b> , <b>Ankh</b>	230M	MLM & CL
PTM-Mamba [125]	UniProt & Swiss-Prot [141] <sup>8</sup>	AA Seq. & PTM Seq.	Mamba [96]	-	MLM

\* **Seq.** - Sequence; **PE** - Positional Encoding; **Encoder** - Transformer Encoder; **Decoder** - Transformer Decoder; **Attn.** - Attention; **Emb.** - Embeddings; **MLM** - Masked Language Modeling; **NTP** - Next Token Prediction; **CRS** - Corrupted Spans Reconstruction; **Recon.** - Reconstruction; **Trans.** - Translation; **CL** - Contrastive Learning;

<sup>1</sup> UniRef database provides clustered sets of protein sequences based on the UniProt Knowledgebase (UniProtKB) and UniProt Archive (UniParc) records. UniRef50, UniRef90, and UniRef100 are clusters of protein sequences at 50%, 90%, and 100% identity. As of December 2024, they have approximately 68 million, 199 million, and 435 million results, respectively.

<sup>2</sup> Big Fantastic Database (BFD) combines UniProt with metagenomic data, containing over 2.6 billion protein sequences. BFD30 clusters the proteins at 30% identity.

<sup>3</sup> ColabFoldDB is established by merging various metagenomic databases and contains over 740 million proteins.

<sup>4</sup> Observed Antibody Space (OAS) collects and annotates over one billion antibody sequences from over 80 studies.

<sup>5</sup> Structural antibody database (SABDab) provides antibody structures that are consistently annotated and presented.

<sup>6</sup> As of December 2024, Pfam database collects approximately 23 thousand protein families, each represented by MSAs and hidden Markov models.

<sup>7</sup> UniClust databases cluster the UniProtKB sequences used on 90%, 50%, and 30% pairwise sequence identity levels.

<sup>8</sup> UniProtKB/Swiss-Prot contains manually reviewed textual annotations of proteins extracted from literature and analysis.

sequence is presented as AA tokens, the protein structure is presented as another track of discrete tokens meanwhile injected into the first transformer block, and then different aspects of protein function (e.g., solvent accessible surface

area, function annotations) are presented as more token tracks. Subsequently, ESM-3 is trained with a special masked language modeling objective: masks are randomly sampled and applied to each track, and the masked tokens should be

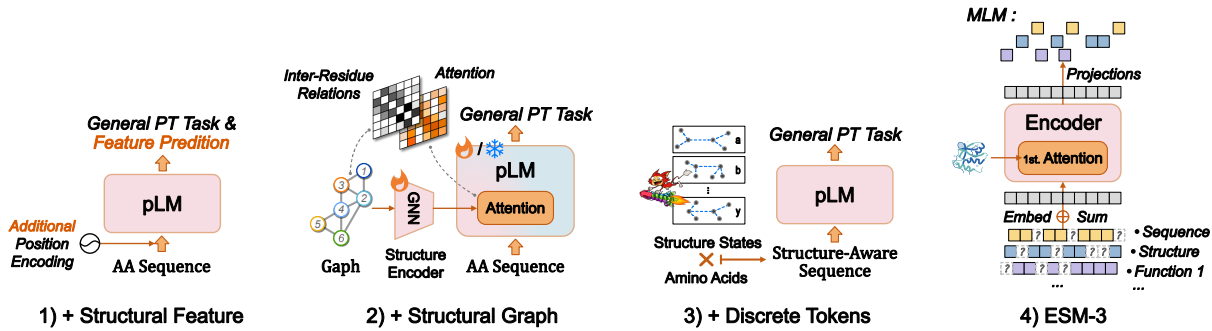


Fig. 5: Typical **structure-enhanced pLMs**. 1) Pre-calculated structural features can be injected into the input AA sequence as position encoding, or utilized in an additional training objective. 2) Considering the significant correlation between the Transformer attention map and protein structural contacts, structural graphs can be encoded by GNNs and combined with the attention module of pLMs. 3) Local structure states along the polypeptide chain are distilled into discrete tokens, which are subsequently involved in the training procedure of pLMs. 4) ESM-3 presents the sequence, structure, and functions of protein as multiple tracks of discrete tokens, with all kinds of information fused within a unified latent space. In particular, there is an additional geometric attention contained in the first Transformer block to process the protein backbone structure.

predicted at the output. Therefore, it learns a single latent space with all kinds of information fused. In evaluations, ESM-3 not only demonstrates excellent benchmarking results but also drives a case of real-world protein design. By using ESM-3 for protein generation, a new bright fluorescent protein is discovered. It is distant from known fluorescent proteins (with  $\sim 58\%$  identity only) and should take 500 million years of evolution to emerge in the wild.

### 3.2.2 Function Enhanced pLMs

The sequence and structure of proteins are clear concepts denoting the linear arrangement and the spatial configuration of residues, respectively. In contrast, protein function represents a multifaceted notion that encompasses diverse categories. As shown in figure 2-F, protein function is labeled by experimental results or manual annotations, and can be documented in academic textual materials.

To empower pLMs with the awareness of functional labels, existing studies guide pLMs to capture: 1) the forward correlation from protein sequences to functional labels, 2) the inverse correlation from functional labels to protein sequences, or 3) the bidirectional inter-correlations between them. Figure 6 illustrates these three categories of methods.

- First, function prediction objectives are introduced within pre-training frameworks. For instance, ProteinBERT [162] combines mask language modeling with GO annotation prediction for pre-training. PromptProtein [163] presents a multi-task pre-training framework tailored for pLMs. The PromptProtein model accepts an AA sequence as input, along with a prompt token that specifies the pre-training task, which can be MLM,  $C_\alpha$  coordinate prediction, or protein-protein interaction prediction. Notably, ESM-1v [72] learns to score the sequence variation effect through a modified masked language modeling objective. With the mutated positions masked, ESM-1v is asked to compare the probability assigned to the mutant with that assigned to the wild-type. After such pre-training, ESM-1v can accurately capture the mutational effects on comprehensive protein function, even in the case of zero-shot inference.
- Second, some pLMs are trained to generate protein sequences conditioned on functional labels. For example,

ProGen [35] enhances 280 million protein sequences from a wide range of families with prefix control tags that specify functional properties, and then proceeds the autoregressive pre-training. As a result, ProGen can generate functional proteins based on the provided control tag across diverse protein families. Similarly, ZymCTRL [164] is a conditional language model specifically pre-trained on the enzyme space, designed to generate enzyme sequences based on user-provided enzyme commission (EC) numbers.

- Third, certain pLMs learn to perform protein function prediction and conditional sequence optimization collaboratively. While taking the concatenated functional property and AA sequence as input, Regression Transformer (RT) [165] undergoes unified pre-training that involves both masked property reconstruction and masked residue span reconstruction. Moreover, ProteinNPT [34] integrates multiple sequence alignment (MSA) and auxiliary functional labels. ProteinNPT is a non-parametric Transformer with tri-axial self-attention across the residue/label tokens, the aligned homologous sequences, and the labeled instances, capable of both functional property prediction and iterative protein redesign.

Text is a flexible and inclusive data modality that can encompass protein functional information. By aligning the representation space of protein sequences with their corresponding textual descriptions, pLMs acquire not only explicit protein function knowledge but also a certain level of natural language understanding. Figure 6-4 illustrates the typical bi-stream autoencoding framework. For instance, Protein-CLAP [166] performs contrastive learning (CL) to align protein-text modalities, thereby producing representations containing protein knowledge from text prompts. These representations could subsequently drive the text-guided protein design. Besides, ProtST [167] is a multi-modal learning framework designed to enhance autoencoding pLMs. During pre-training, pLMs learn through the incorporated unimodal masked language modeling (MLM), cross-modal MLM, and protein-text CL. It is revealed that ProtST-induced pLMs outperform the vanilla pLMs in representation learning benchmarks.

Moreover, as one of the most successful contrastive

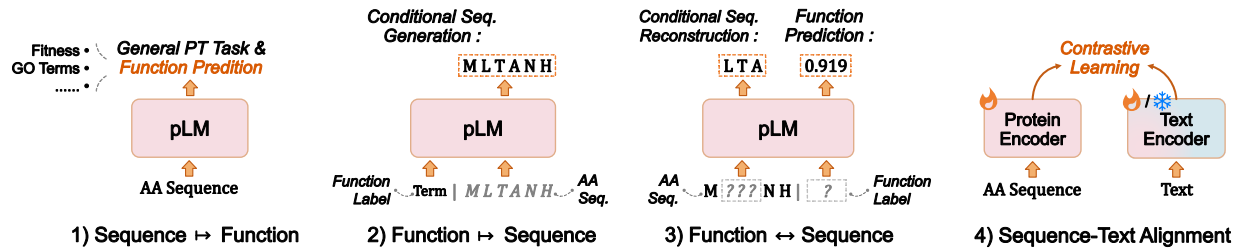


Fig. 6: Typical **function-enhanced pLMs**. 1-3) pLMs learn the forward, inverse, and bidirectional correlations between protein sequences and functional labels. They undergo pre-training using various objectives, including function prediction, conditional sequence generation, or a variant of mask language modeling that combines the two. 4) Protein sequences are aligned with their corresponding textual descriptions through contrastive learning.

learning methods, CLIP is first proposed to connect the text and images [86], training a text encoder and an image encoder to predict the pairings of a batch of  $(text, image)$  examples. In computational protein science, the technical framework of CLIP is employed to distill the same protein semantics described in different vocabularies. Based on the curated pairs of AA sequence and function annotation text, ProtET [168] performs CLIP-like pre-training to align features of protein sequence and text, and ProteinCLIP [169] equips adapter layers for the pre-trained pLM and text encoder that each re-project the protein or text representations into a shared representation space. Furthermore, ProTrek [170] is a tri-modal pLM that conducts contrastive learning of the sequence, structure, and function text. In addition to the comprehensive representation of proteins, ProTrek supports all kinds of cross-and intra-modal retrieval with a total of nine searching tasks. We are allowed to precisely navigate the vast protein universe in seconds by giving the protein sequence, structure, or even natural language description.

Besides, knowledge graphs (KGs) provide factual knowledge of protein functions that further integrate classification annotations and textual documentation. ProteinKG25 [171] is a KG dataset that includes protein entities with sequences, Gene Ontology (GO) term entities with textual descriptions, as well as the GO-GO and protein-GO triples. On this basis, OntoProtein [171] performs incremental pre-training for ProtBERT through MLM and a knowledge embedding objective, and KeAP [172] performs protein-knowledge exploration at the token level with well-designed cross-attention modules. To compare those recent advances, Ko et al. [173] present a benchmark for functional text-integrated pLMs. In assessing the learned representations, the authors implement six models (ProteinCLAP, ProtST, ProteinCLIP, ProTrek, OntoProtein, and ESM-3) with a sequence-only baseline pLM (ESM-2) across six downstream tasks. It is revealed that these function-enhanced pLMs outperform ESM-2 in five of six tasks, while no one is always the best.

Lastly, we simply review some relevant protein representation learning studies that aim for a comprehensive fusion of protein semantics. In integrating sequence and structure, SSEmb [174] combines a pLM for multiple sequence alignments with a graph representation for the protein structure. ESM-GearNet [59] explores different fusion strategies, including serial, parallel, and cross fusion, to combine sequence representations from pLMs (such as ESM-2) with structure representations from structure encoders

(like GearNet). Lee et al. [175] introduce a pre-training strategy that incorporates multi-view protein knowledge from the sequence, 3D structure, and surface for improved representation learning. BioCLIP [176] is a contrastive learning framework designed to train protein structure models by utilizing pLMs for assistance. Pursuing the unity of protein sequence-structure-function, MASSA [177] aligns independently derived embeddings of AA sequence, structural graph, and GO terms. Additionally, Protein-Vec [178] is a multi-view information retrieval system for proteins, where a mixture-of-experts model is employed to combine protein sequences with seven structural and functional properties.

### 3.3 Multimodal pLMs

In the preceding subsections, we have introduced existing pLMs that decipher protein sequences and understand the structural and functional information. While some of these models incorporate textual descriptions linked to proteins, their primary focus remains on protein-centered semantics. In this subsection, we introduce pLMs that exhibit proficiency in external languages, encompassing natural language with world knowledge, chemical molecule language, and beyond. As those languages convey greatly varied semantics, here we recognize them as different modalities. Those multimodal pLMs are summarized in Table 3. Meanwhile, Figure 7 illustrates two typical technical approaches: unified multimodal learning, or "specific encoder - unified decoder" models.

**Protein-Text Models** present the emergent ability of general LLMs and exhibit their superiority in instruction understanding and multi-task processing. InstructProtein [188] undergoes pre-training on both natural language and protein sequence corpora, followed by instruction tuning to align between these two distinct languages. ProtLLM [189] designs a protein-as-word language modeling approach, associating named protein entities in broad biological corpus with their AA sequences. This approach effectively unifies the protein and word inference as an autoregressive next-token prediction task. ProLLaMA [31] is a novel training framework that empowers general LLMs with the capability of protein language processing. Based on the low-rank adaptation (LoRA) [82] technique, this framework consists of two stages: incremental training using protein sequences and instruction tuning on the instruction dataset. Furthermore, Evolla [190], a multimodal pLM with 80 billion parameters, is developed for unraveling protein molecular mechanisms via natural language dialogue. To overcome the quantity and



**TABLE 2: Structure and Function Enhanced pLMs.** For each pLM, the resource link, pre-training corpora, input data format, network architecture, size of full parameters, and pre-training objective is summarized here. Unusual abbreviations are explained in the footnote \*, and important databases are briefly introduced in footnotes 1 to 9. In the "Architecture" field, models employed in a frozen form are colored **blue**, models fine-tuned are colored **pink**, and the others are trained from scratch.

Method	Corpora	Input	Architecture	#Parameter	Objective
PeTriBERT [142]	AlphaFoldDB [179] <sup>1</sup>	AA Seq. + Struct. PE	BERT	40M	MLM
ESM-S [143]	DeepSF [180]	AA Seq.	<b>ESM-2</b>	8M ~ 650M	Remote Homology Detection
METL [144]	-	AA Seq. + Struct. PE	Encoder	20M, 50M	Biophysics Pred.
LM-Design [148]	CATH [181] <sup>3</sup>	AA Seq. & Struct. Graph	<b>ESM-1b</b> w/ Struct. Adapter	658M	Struct.-Cond. MLM
ProseLM [149]	CATH / PDB [45] <sup>2</sup>	AA Seq. & Struct. Graph	<b>ProGen2</b> w/ Struct. Adapter	151M ~ 6.4B	Struct.-Cond. NTP
PST [147]	AlphaFoldDB	AA Seq. & Struct. Graph	<b>ESM-2</b> w/ Struct. Extractors	8M ~ 650M	MLM
ProstT5 [155]	AlphaFoldDB	AA Seq. & 3Di Tokens	<b>ProtT5</b>	3B	CSR & AA-3Di Trans.
SaProt [29]	AlphaFoldDB & PDB	Struct.-Aware Tokens	ESM-2	35M, 650M	MLM
ProSST [157]	AlphaFoldDB	AA Seq. & Struct. Tokens	Encoder w/ Disentangled Attn.	110M	MLM
DPLM-2 [161]	PDB & AlphaFoldDB	AA Seq. & Struct. Tokens	<b>DPLM</b>	150M ~ 3B	Discrete Diffusion
ESM-3 [30]	UniRef, MGnify [126], JGI [127], OAS, PDB, AlphaFoldDB, ESMAtlas [25] <sup>4</sup> , InterPro [182] <sup>5</sup> , InterProScan [183] <sup>5</sup>	AA Seq., Struct. Tokens & Func. Tokens	Encoder w/ Struct. Block	1.4B ~ 98B	MLM
ProteinBERT [162]	UniRef90 & GO [46] <sup>6</sup>	AA Seq. & GO Annot.	Customized Encoder	16M	NPT & GO Pred.
PromptProtein [163]	UniRef50, PDB & STRING [184] <sup>7</sup>	AA Seq.   Prompt Token	Encoder w/ Prompt-Aware Attn.	650M	MLM, Struct. Pred. & PPI Pred.
ESM-1v [72]	UniRef90	AA Seq.	ESM-1b	650M	Mutation Scoring MLM
ProGen [35]	Pfam, GO, & NCBI Taxonomy [185]	Func. Tag   AA Seq.	Decoder	1.2B	Func.-Cond. NTP
ZymCTRL [164]	BRENDA <sup>8</sup>	EC Number   AA Seq.	Customized Decoder	738M	Func.-Cond. NTP
RT [165]	TAPE [158]	Prop. Label   AA Seq.	XLNet w/ Bidirec. Attn.	27M	CSR & Prop. Pred.
ProteinNPT [34]	-	MSA   Prop. Label	Encoder w/ Tied Row & Column Attn.	-	MLM & Prop. Pred.
ProteinCLAP [166]	Swiss-Prot [141] <sup>9</sup>	AA Seq. & Func. Text	<b>ProtBERT</b> , <b>SciBERT</b> [186]	420M	Seq.-Text CL
ProtET [168]	Swiss-Prot & TrEMBL [141] <sup>9</sup>	AA Seq. & Func. Text	<b>ESM-2</b> , <b>PubMedBERT</b> [187]	750M	Seq.-Text CL
ProteinCLIP [169]	UniProtKB	AA Seq. & Func. Text	<b>ESM-2</b> / <b>ProtT5</b>	8M ~ 3B	Seq.-Text CL
ProTrek [170]	UniRef50, AlphaFoldDB, PDB, & Swiss-Prot	AA Seq., Struct. Tokens & Func. Text	<b>ESM-2</b> , Encoder, <b>PubMedBERT</b>	930M	Seq.-Struct.-Func. CL
ProtST [167]	Swiss-Prot	AA Seq. & Func. Text	<b>ProtBERT</b> / <b>ESM-1b</b> / <b>ESM-2</b> , <b>PubMedBERT</b>	420M / 650M	MLM & CL
OntoProtein [171]	ProteinKG25 [171]	AA Seq. & KG	<b>ProtBERT</b> , <b>PubMedBERT</b>	420M	KE & MLM
KeAP [172]	ProteinKG25	AA Seq. & KG	Encoder-Decoder, <b>PubMedBERT</b>	-	MLM

\* **Struct.** - Structure; **Func.** - Function; **Prop.** - Functional Property; **Annot.** - Annotation; **KG** - Knowledge Graph; **Attn.** - Attention; **Pred.** - Prediction; **Cond.** - Conditioned; **Trans.** - Translation; **PPI** - Protein-Protein Interaction; **KE** - Knowledge Embedding

<sup>1</sup> AlphaFold Database (AlphaFoldDB) provides reliable structure predictions for over 200 million proteins.

<sup>2</sup> Protein Data Bank (PDB) collects more than 220 thousand experimentally-determined 3D structures of proteins.

<sup>3</sup> CATH is a classification of protein structures based on PDB. 151 million protein domains are grouped into over 5 thousand superfamilies.

<sup>4</sup> ESM Metagenomic Atlas (ESMAtlas) contains 772 million metagenomic protein structures predicted by ESMFold.

<sup>5</sup> InterPro and InterProScan provide functional analysis of proteins by categorizing them into families and predicting their domains and crucial sites.

<sup>6</sup> Gene Ontology (GO) knowledgebase systematically annotates protein functions with GO terms.

<sup>7</sup> STRING integrates all public protein-protein interaction (PPI) sources, covering over 2 billion PPIs among 59.3 million proteins.

<sup>8</sup> BRENDA database contains approximately 37 million enzyme sequences and their corresponding EC number annotations.

<sup>9</sup> UniProtKB/Swiss-Prot provides manually reviewed protein annotations; UniProtKB/TrEMBL consists of proteins with computationally analyzed annotations.

quality bottlenecks in the paired protein-text data, Evolla is trained on an unprecedented AI-generated dataset involving 546 million protein-question-answer triples, which could narrow the scale gap between protein sequences and the corresponding functional annotations.

In biological systems, protein-molecule interactions play a pivotal role, which is especially important for medicine. This emphasizes the significance of bi-lingual **Protein-Molecule Models**. DrugGPT [191] is an autoregressive LM trained on a substantial amount of protein-ligand binding data. It can

generate potential ligand SMILES conditions corresponding to specific protein sequences. ESM All-Atom (ESM-AA) [192] enables unified molecular modeling at both the residue scale for proteins and the atom scale for small molecules. In the development of ESM-AA, the main technical contributions include a multi-scale position encoding scheme that captures relationships among residues and atoms, and a pre-training framework for code-switching-enhanced protein sequences.

Furthermore, **Protein-Text-Molecule Models** aim to decipher multiple biomolecular languages in a consistent

TABLE 3: **Multimodal pLMs**. In this table, we present the resource link, pre-training corpora, input data format, network architecture, scale of parameters, and learning procedure for typical *protein-text*, *protein-molecule*, *protein-text-molecule*, and *broader modal* models. In the context of "Corpora", important databases are briefly described in footnotes 1 to 5. In the "Architecture" field, parameters in uncolored models are trained from scratch, pre-trained parameters in pink models are trainable, and pre-trained parameters in pink models are frozen.

Method	Corpora	Input	Architecture	#Parameter	Learning Procedure
ProtLLM [189]	UniProt, PubMed [193] <sup>1</sup> , STRING & Mol-Instructions [194]	AA Seq. & Text	ProtST, LLaMA-2 [18]	7B	Protein-as-word Pre-training
ProLLaMA [31]	UniRef50 & InterPro	AA Seq. & Text	LLaMA-2	7B	Instruction Tuning
InstructProtein [188]	UniRef100 & PubMed	AA Seq. & Text	Decoder	1.3B	Multimodal Pre-training & Instruction Tuning
Evolla [190]	Swiss-Prot, ProTrek [170], TrEMBL & AI-generated data	AA Seq. & Text	SaProt, LLaMA-3 [195], Align Module	10B, 80B	Causal Language Modeling & Direct Preference Optimization
DrugGPT [191]	ZINC20 [196] <sup>2</sup> & Jglaser <sup>3</sup>	AA Seq. & SMILES	GPT-2	1.5B	Ligand & Protein-Ligand Pre-training
ESM-AA [192]	AlphaFoldDB & Uni-Mol [197]	AA Seq. & SMILES + Multi-scale PE	ESM-2	35M	Pre-training w/ MLM & Pairwise distance recovery
BioMedGPT [198]	UniProt & PubChem [199] <sup>4</sup>	AA Seq., Text & Mol. Graph	GraphMVP [200], ESM-2, LLaMA-2	10B	Multimodal Fine-tuning
BioT5 [32]	UniRef50, ZINC20, PubChem, PubMed, C4 [14] <sup>5</sup> & Swiss-Prot	AA Seq., Text & SELFIES	T5	252M	Pre-training w/ CSR & Multimodal Translation
BioT5+ [201]	UniRef50, ZINC20, PubChem, PubMed, C4 & Swiss-Prot	AA Seq., SELFIES, IUPAC & Text	T5	252M	Pre-training w/ CSR & Multimodal Translation
InstructBioMol [202]	PubMed, bioRxiv, ChemRxiv, PubChem, UniRef50, ChEBI-20 [203], TrEMBL, Swiss-Prot, BindingDB [204] & Rhea [205]	SELFIES, Mol Graph, AA Seq., Struct. Graph & Text	Struct. GNNs, ESM-2, SaProt, LLaMA-2	6.8B	Continual Pre-training & Instruction Tuning
BioTranslator [206]	STRING, ChEBI-20 & Human Phenotype Ontology [207]	Text, AA Seq., Gene Expression & Phenotype Pathway	PubMedBERT, Non-Text Projectors	100M	Contrastive Learning
Galactica [208]	-	AA Seq., SMILES, DNA Seq., & Code	Decoder	120B	Prompt Pre-training

<sup>1</sup> PubMed database contains the abstract and citation information for more than 37 million biomedical literature.

<sup>2</sup> ZINC20 is a chemical database that contains billions of molecules and supports precious searching.

<sup>3</sup> Jglaser dataset contains 1.9 million unique pairs of protein sequence & ligand SMILES with experimentally determined binding affinities.

<sup>4</sup> PubChem provides information about chemical molecules, such as the SMILES string and IUPAC names.

<sup>5</sup> Colossal Clean Crawled Corpus (C4) is a collection of general English-language text sourced from the public web scrape.

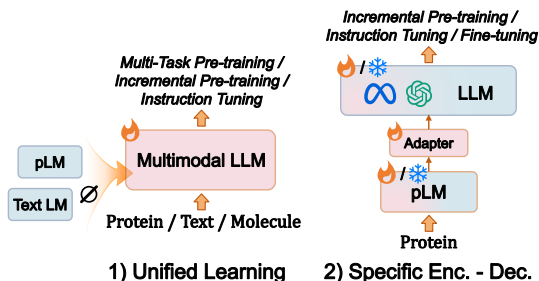


Fig. 7: Typical **multimodal pLMs**. 1) In unified multimodal learning, multiple scientific languages share a unified latent space. We can perform multi-task pre-training from scratch, as well as incremental pre-training or instruction tuning based on pre-trained LMs. 2) In "specific encoder - unified decoder" models, one or more specific encoders are connected to a unified decoder, with internal adapters typically learning to bridge the semantic spaces. The end-to-end model can learn through incremental pre-training, instruction tuning, or fine-tuning.

manner. BioMedGPT [198] aligns the feature spaces of an autoencoding pLM, a molecular graph encoder, and an autoregressive LLM by using a question-answering-based

fine-tuning approach. BioT5 [32] is a pre-training framework proposed for comprehensive multi-language integration. The BioT5 model is trained across data of protein sequences, molecule SELFIES strings, scientific texts, and wrapped sentences using the corrupted span reconstruction and bidirectional translation objectives. Then, BioT5+ [201] incorporates IUPAC names to enhance molecular understanding, aiming to bridge the gap between specialized molecular representations and the corresponding textual descriptions. Notably, Mol-Instructions [194] is a meticulously curated dataset that consists of protein-oriented, molecule-oriented, and biomolecular text instructions. This dataset facilitates the adaptation of general LLMs, such as the LLaMA series, by enabling effective instruction tuning to address diverse biomolecular tasks. Moreover, InstructBioMol [202] is a novel LLM learned through a comprehensive any-to-any alignment of natural language, molecules, and proteins. InstructBioMol has exhibited critical capabilities of biomolecular instruction following and multimodal data understanding, and demonstrated the potential to serve as a digital research assistant in supporting practical biomolecular tasks.

Besides, there have been extensive explorations for broader scientific languages, as well as investigations into

the methods for bridging multiple languages. For instance, BioTranslator [206] learns a cross-modal translation to bridge the user-written text and the corresponding non-text biological data, such as protein sequences, gene expression vectors, and phenotype pathways. Galactica [208] is an LLM capable of storing, combining, and reasoning about scientific knowledge. It demonstrates proficiency in handling diverse data modalities, including general text, LATEX code, programming code, molecular SMILES, protein sequences, and DNA sequences. Notably, BioBridge [209] is a parameter-efficient learning framework that bridges independently pre-trained scientific LLMs based on the biomedical knowledge graph. In constructing multimodal scientific LLMs, BioBridge presents promising to overcome the limitations imposed by data scarcity and computational costs.

## 4 UTILIZATION AND ADAPTATION OF PROTEIN LANGUAGE MODELS

While the presence of progressive information flow between protein sequence, structure, and function is widely acknowledged, the intricate principles underlying protein folding and functioning are still awaiting revelation. Structure prediction, function prediction, and protein design have long been fundamental modeling problems in the field of computational protein science. As pLMs possess a fundamental understanding of proteins, they have been utilized or adapted to these problems, yielding significant advancements and impacts.

### 4.1 Protein Structure Prediction

In structural biology, scientists determine protein structures through experimental techniques like X-ray crystallography [41], nuclear magnetic resonance [210], and electron cryomicroscopy [42]. These laboratory works are complex, time-consuming, and expensive. It is estimated that determining each protein structure takes months to years and costs tens of thousands of dollars [211]. So far, there are only about two hundred thousand experimentally determined structures collected in the Protein Data Bank [45]. At this rate of development, it would take millions of research-years to analyze hundreds of millions of natural proteins that are sequenced but of unknown structures. In this context, protein structure prediction emerges as a crucial challenge. If a computational model can *accurately infer the atom-wise 3D structures of proteins from their amino-acid sequences*, the progress of human understanding of protein structures would be significantly accelerated.

In recent years, rapid developments in artificial intelligence and computing power have greatly promoted the advancement of protein structure prediction. Breakthrough methods like AlphaFold2 [8] and RoseTTAFold [9] exhibit an unprecedented level of near-experimental accuracy in predicting protein structures. They have played the role of essential tools for scientists to obtain a reliable protein structure within tens of minutes. The great success of these models should be credited to the ingenious incorporation of network architectures and training strategies, which well meet the evolutionary and geometric constraints on protein structures. Elaborate workflows are proposed to extract co-evolution knowledge from multiple sequence alignments

(MSAs) and enforce the pairwise description of residues to satisfy the triangle inequality on distances.

For instance, Figure 8-1 illustrates the workflow of AlphaFold2. Given an input amino-acid sequence, AlphaFold2 first retrieves an MSA and a small number of homologous structures (i.e., templates). Notably, experience suggests that the MSA would significantly impact the model's performance more than the templates do. After initial encoding, the inference begins with an MSA representation that conveys evolutionary constraints and a pairwise representation that embodies geometric information. In the following Evoformer blocks, the two representations undergo constraint transformations and are updated mutually, which enables sufficient reasoning and fusion about the evolutionary and geometric information. Subsequently, all-atom positions and side-chain angles are inferred by a structure module. Besides, the model undergoes "recycling" iterations: the training loss is applied to outputs repeatedly, while the outputs are input back into the same modules recursively. Overall, within the latent space of AlphaFold2, the 3D structure of a protein is first hypothesized by an MSA and then progressively refined into an increasingly accurate prediction result.

Although these methods have demonstrated their effectiveness in protein structure prediction, they still encounter challenges due to the reliance on MSAs. First, searching for MSA from extensive databases is time-consuming, usually taking up most of the total inference time. It should be certain that speeding up protein structure prediction would further broaden its applications. Second, it's inherently difficult to obtain sufficient MSAs for orphan or fast-evolving proteins that lack homology. The reliance on MSA impedes the accuracy of structure prediction for these special proteins. Third, from a theoretical point of view, protein folding is an independent process for each protein. All structural information of a protein is fully encoded in its single sequence rather than the retrieved MSA. Models like AlphaFold2 learn MSA-to-structure but not sequence-to-structure, which deviates from the ultimate goal of grasping the rule of protein folding.

To overcome these MSA-caused challenges, methods for single-sequence (MSA-free) protein structure prediction have been developed, with pLMs playing a foundational role. Through large-scale pre-training, pLMs have learned the favorable patterns of protein sequence and some implicit knowledge of protein structure. When given a single amino-acid sequence as input, pLMs encode evolutionary constraints in sequence representations and embody the inter-residue structural relationships in attention maps (i.e., pairwise representations). In conveying the same information, the traditional MSA can be replaced by the sequence and pairwise representations generated by pLMs. This supports subsequent deep-learning modules to model protein folding as well. Thus, from an end-to-end view, the typical workflow of single-sequence protein structure prediction involves taking pLMs as the foundation model and training additional "prediction heads" to deduce geometric constraints and complete structure predictions.

A series of methods build AlphaFold2-like "prediction heads" comprising a geometric module (i.e., EvoFormer) and a structure module, with the geometric module as the focus of improvement. For example, the workflow of ESMFold [25] is presented in Figure 8-2. During inference,

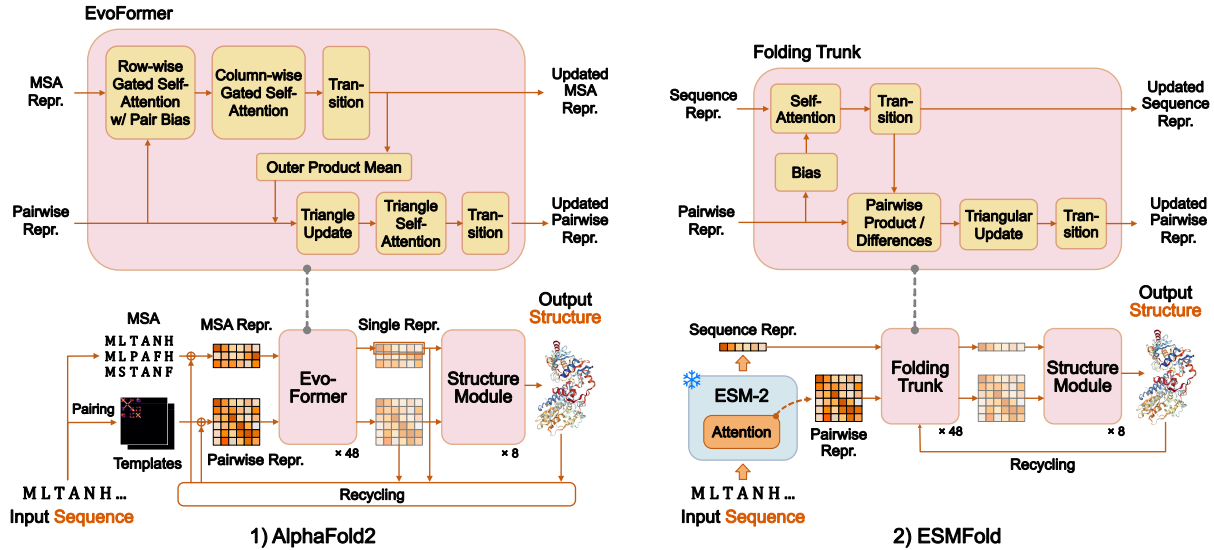


Fig. 8: **Workflow overview of AlphaFold2 [8] and ESMFold [25].** Both AlphaFold2 and ESMFold infer the high-resolution protein structure from a comprehensive understanding of the sequence. AlphaFold2 relies on MSA to gain evolutionary insights encoded in protein sequences, while ESMFold achieves this through the utilization of a protein language model.

the sequence of a protein is fed into the foundation model (ESM-2), and the internal representations are extracted and sent to the folding trunk that contains an array of folding blocks. Resembling the EvoFormer block of AlphaFold2, each ESMFold's folding block updates the sequence representation and the pairwise representation alternately. Outputs of the folding trunk are then directed to an equivariant structure module, undergoing three rounds of recycling until the final protein structure prediction is produced. So far, ESMFold has been applied to predict structures for more than 617 million protein sequences, which provides a fire-new view into the diversity of natural proteins at the evolutionary scale. Furthermore, xT-Fold [27] has a similar architecture to ESMFold, while the number of stacked folding blocks is reduced from 48 to 1. HelixFold-Single [62] revises the Evoformer of AlphaFold2 into an EvoformerS (Evoformer with Single representations) module, where the column-wise gated self-attention originally designed for MSA is removed. OmegaFold [212] introduces Geoformer, a novel transformer neural network inspired by geometry principles. In Geoformer layers, the sequence and pairwise representations are iteratively smoothed, and the geometric inconsistency among them is gradually reduced.

In addition, certain methods design distinct "prediction heads" that involve geometric modeling beyond updating the sequence and pairwise representations alternatively. In trRosettaX-Single [61], the sequence and pairwise representations are concatenated and passed to Res2Net-Single, a multi-scale neural network that distills inter-residue 2D geometry. The predicted inter-residue distance and orientations are then utilized to reconstruct 3D structures through energy minimization. In RGN2 [213], a pLM is combined with a Recurrent Geometric Network that uses Frenet-Serret frames to generate protein backbone structures. Furthermore, IgFold [214] is a specific antibody structure prediction model. IgFold consists of an antibody LM (i.e., AntiBERTy), followed by a series of customized graph networks and transformer

layers that predict the backbone atom coordinates directly.

Compared to mainstream MSA-based methods, these pLM-based single-sequence structure prediction methods exhibit the advantages of being fast and universal. On the one hand, these methods are compared to AlphaFold2 and RoseTTAFold in terms of the inference runtime. For big proteins made up of several hundreds of amino acids, single-sequence protein structure prediction often cuts the runtime to less than one-third. While for short sequences of dozens of amino acids, there can be a hundredfold speed advantage. On the other hand, the structure prediction performance for proteins with no homology is also improved. In investigations, OmegaFold, trRosettaX-Single, and RGN2 are all observed to outperform AlphaFold2 and RoseTTAFold on those orphan proteins and *de novo* designed proteins. However, there still remains room for improvement in accuracy: single-sequence protein structure prediction achieves competitive accuracy but is not superior to the MSA-to-structure prediction. Such a result indicates that the principle of protein folding is still not sufficiently grasped.

Recently, AlphaFold3 [215] has been proposed as the latest updation of the AlphaFold series, making surges in the field of protein science again. In addition to predicting protein structures solely, AlphaFold3 can infer the joint 3D structure of complexes, covering nucleic acids, small molecules, ions, and modified residues, and has demonstrated substantially improved accuracy over previous state-of-the-art specialized docking or interaction prediction tools. Technically, the advancement is achieved by the simpler but more general network architecture and training procedure. The reliance on MSA is significantly reduced by replacing the Evoformer of AlphaFold2 with a novel Pariformer module, where the MSA representation is removed, and all information is passed through the pairwise representation for geometric modeling. Moreover, the structure module of AlphaFold2 is replaced by a diffusion module to predict the raw atomic coordinates directly. Compared to AlphaFold3, pLM-based single-sequence protein structure prediction is



advancing in a relatively consistent direction to stride the traditional MSA processing. However, there is still a long way to go toward understanding the fundamental physical and chemical laws of molecules.

## 4.2 Protein Function Prediction

Unlike the clearly defined protein sequence and structure, protein function exhibits multifaceted characteristics as different proteins play diverse biological roles in broad living systems. Then, greater gaps exist between the number of sequenced or structure-determined proteins and ground-truth function labels. In the Swiss-Prot [216] database, less than 100 thousand proteins are manually annotated with Gene Ontology (GO) terms. For specific functional properties, such as stability, fluorescence, and fitness values, typical benchmark datasets (e.g., TAPE [158], FLIP [63], PEER [159], ProteinGym [160]) consist of roughly tens of thousands of labeled instances. Certain instances are under extreme data scarcity as the wet lab experiments have only been conducted on a minimal scale. As a result, protein function prediction has presented significant value in guiding the study of proteins that lack corresponding functional knowledge, where a wide range of different prediction tasks are involved.

Before the emergence of pLMs, AI models are individually trained from scratch for various protein function prediction tasks. This traditional paradigm had a tough drawback: as the models lack transferrable protein knowledge, the prediction performances are barely satisfactory, especially in the case of data scarcity. To overcome the shortcoming, pLMs have been successfully utilized or adopted in protein function prediction. Figure 9 illustrates the typical utilization or adaptation techniques. In the "*LM-as-Encoder*" scheme, prediction heads are typically incorporated following the main body of autoencoding pLMs. These heads are either trained independently with pLM parameters frozen or fine-tuned alongside the pLM in full-model or parameter-efficient manners. Even more simply, the likelihood probability derived in language modeling is skillfully transferred into predictions. In the "*LM-as-Predictor*" scheme, a unified LLM generates answers according to prompts specifying contexts and questions. Consequently, the acquired knowledge within pLMs is effectively transferred to downstream predictions, resulting in significant advancements across various specific protein function prediction tasks. We organize this subsection by summarizing these tasks into five categories: functional property prediction, functional class annotation, functional site identification, protein-molecule interaction prediction, and multi-task question answering. Meanwhile, Table 4 presents a comprehensive summary of relevant studies.

### 4.2.1 Functional Property Prediction

In **functional property prediction**, each protein  $x$  is mapped to a label  $y \in \mathbb{R}$  that measures a quantitative functional property, and existing studies generally employ the *LM-as-encoder* scheme. For example, LMDisorder [217] utilizes the representations generated by ProtT5 to predict the disorder probability of proteins. On the basis of fine-tuning ESM-2, ESMtherm [218] predicts the folding stability of proteins, PIC [220] predicts human protein essentiality, VISH-Pred [219] is an ensemble framework for protein toxicity

prediction. Specifically, PeptideBERT [221] fine-tunes ProtBERT to predict three critical properties of peptides, i.e., solubility, hemolysis, and non-fouling. LassoESM [222] is a tailored pLM for enhanced lasso peptide property prediction, where the properties include lasso cyclase substrate tolerance, RNA polymerase inhibition activity, etc. Moreover, as the pH level of the reaction environment could significantly affect enzyme activity, OphPred [223] and EpHod [224] are recently proposed to predict the enzyme optimal pH.

Notably, the fitness of a protein is a synthetic property of how well a protein can perform its function within an organism [160]. Experimentally, an assay of deep mutational scanning (DMS) [43] measures the effects of hundreds to thousands of mutants on a single protein. DMS assays are performed for different proteins with varied functions using various forms of experimental measurements [72]. Relevant influencing factors include stability, folding efficiency, binding affinity, etc. That is to say, instead of describing a specific property presented in all proteins, fitness labels assess the comprehensive function performance of a limited set of homologous proteins. Meanwhile, the effect of mutants indicates how the function of mutants changes compared to the wild-type, which is tightly connected to fitness.

Recent studies have revealed that pLMs play a significant role in the zero- and few-shot predictions of protein fitness and mutant effects. After pre-training on massive protein sequences that survived through natural evolution, pLMs understand what sequences are plausible (i.e., probably encode proper function) while the others are invalid. Practical results demonstrate that *the likelihoods inferred from pLMs correlate well with protein fitness* [72, 253, 254]. By comparing the language modeling probabilities assigned to the mutant sequence and wild-type sequence, ProteinGym [160] calculates the mutant effect on fitness based on existing pLMs in a zero-shot manner, and Brandes et al. [225] have predicted all possible missense variant effects in the human genome with an ESM-1b-based workflow.

Besides, fine-tuning pLMs on certain labeled protein fitness values leads to more robust predictions. Few-Shot Learning for Protein Fitness Prediction (FSFP) [226] is an effective training strategy designed to optimize pLMs under extreme data scarcity. By combining advanced techniques of meta-transfer learning, learning to rank, and LoRA, FSFP significantly boosts the performance of pLMs in predicting protein fitness based on only tens of labeled single-site mutants. Then, Denoising Protein Language Models (DePLM) [227] predicts mutation effects by refining evolutionary information captured in pLMs. Conceptually, evolutionary information could comprise both property-relevant and irrelevant information, with the latter acting as "noise" for the specific prediction task at hand. DePLM contains an ingenious denoising diffusion framework for the likelihoods produced by pLMs, effectively filtering out the irrelevant information to improve mutation effect predictions.

### 4.2.2 Functional Class Annotation

**Functional class annotation** can be acknowledged as multi-class, hierarchical, or multi-label classification problems in different situations. In annotating prokaryotic viral proteins, each viral sequence is mapped to one of nine key functional classes. Flamholz et al. [228] construct a pLM-based classifier,

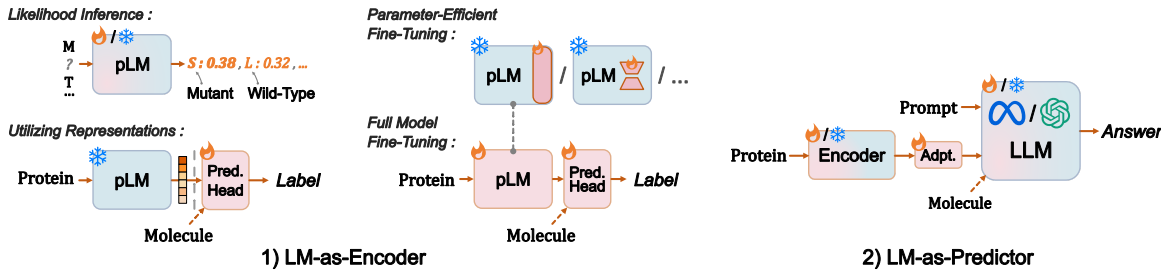


Fig. 9: Typical technical schemes in pLM-based protein function prediction. 1) In the "LM-as-Encoder" scheme, protein language models play a central role in the full models as encoders. Predictions could be inferred from the likelihood of language modeling or obtained through additional prediction heads. 2) In the "LM-as-Encoder" scheme, textual LLMs play a central role in the full models as predictors. In addition to the textual prompts, LLMs receive the encoded protein sequence/structure or molecule knowledge.

enabling fresh regions within the viral sequence space to be annotated meaningfully. In annotating enzymes, each enzyme sequence is mapped to a unique EC number. Fernandez et al. [64] explore different representation strategies of manual feature engineering, frequency space encoding, and pLM encoding, and compare various alternative classifiers in the prediction head. Besides, GraphEC [229] performs EC number prediction based on geometric graph learning, where the featured graphs are constructed using the pLM-produced sequence representations and ESMFold-predicted structures.

In Gene Ontology (GO) function prediction, each protein  $x$  is mapped to a set of labels  $\{y_1, y_2, \dots, y_n\}$  where  $y_i \in \{0, 1\}$  denotes whether the  $i$ -th GO term is annotated to the protein. In recent years, the NetGO series [255, 256, 230] has exhibited remarkable performance in large-scale functional annotations, where pLMs have driven the latest improvement. On the basis of NetGO 2.0, the NetGO 3.0 [230] incorporates a new model, LR-ESM, which leverages ESM-1b to present each protein and undergoes training with logistic regression. Similarly, SPROF-GO [232] leverages the pLM-based representation for initial prediction, and that is further advanced through a label diffusion algorithm. Specifically, there are hierarchical inclusion relationships between GO terms, i.e., the semantic proximity between multi-class labels should be considered in the prediction. Based on representations produced by ESM-2, DeepGO-SE [231] generates multiple approximate models of GO, which predict the truth values of statements about protein functions. The truth values over those models are aggregated to ultimately determine GO functions.

#### 4.2.3 Functional Site Identification

On a micro level, proteins perform functions by some significant sites rather than the whole macromolecule. Given a protein, **functional site identification** aims to reveal which residues are of functional sites, and simultaneously identify the specific type of function if necessary. For example, GraphBepi [233] is a graph model that predicts binding probabilities of B-cell epitope for each residue. Starting with a protein sequence, a protein graph is constructed based on the AlphaFold2-predicted structure, where pLM-learned per-token sequence representations are specifically taken as node attributes. A deep learning model consisting of GNNs and Bi-LSTM is subsequently trained for the prediction. Then GPSite [66] is a multi-task network aiming to predict binding residues of biomolecules (e.g., DNA, RNA, peptide,

etc.) on proteins. The predicted structure from ESMFold and the representation produced by ProtT5 are similarly processed into a geometric-aware attributed protein graph. Subsequently, a shared GNN captures the common binding-relevant characteristics, and individual prediction heads are trained for ligand-specific binding site prediction.

Besides, signal peptides (SPs) are short sequences that regulate protein secretion and translocation across all organisms, currently categorized into five types. SignalP 6.0 [234] combines a pLM with a conditional random field model to identify which residues belong to the SP region and simultaneously infer the SP type. Then, Unbiased Organism-agnostic Signal Peptide Network (USPNet) [235] is another deep learning model proposed for SP classification and cleavage-site prediction. In USPNet, a pLM like ESM-MSA-1b or ESM-2 serves as an encoder to enrich representations of sequences, which facilitates the prediction derived by additional learning modules.

#### 4.2.4 Protein-Molecule Interaction Prediction

Proteins generally carry out functions in concert with additional molecules, which can be another protein, DNA, RNA, or drug. In addition to identifying the binding sites solely on proteins, predicting the protein-molecule interaction states is equally important. Given paired protein and molecule as input, **protein-molecule interaction prediction** is designed to determine whether they interact, measure relevant properties about the interaction, or precisely identify the interaction sites. For example, ConPLex [67] is proposed to predict the interaction probability between drugs and protein targets. ConPLex aligns the embedded molecular fingerprints with the pLM-generated protein representations by contrastive learning, and calculates the distance between a pair of aligned representations for binary prediction. Similarly, using language models to encode proteins and ligands, BALM [238] learns protein-ligand binding affinities by optimizing their cosine similarity in a shared latent space. Then, while taking pLM-produced representations and molecular graphs as inputs, BIND [236] can predict the drug-target affinity values and discriminate between active and decoy ligands.

Notably, enzymes are workhorses for various biological processes, binding to substrates and releasing products in diverse catalyzed reactions. In understanding the functionality of enzymes, UniKP [68] aims to predict enzyme kinetic parameters (i.e., enzyme turnover number, Michaelis constant, catalytic efficiency) from protein sequences and

TABLE 4: **pLM-based Protein Function Prediction Methods.** Corresponding to the technical schemes summarized in Figure 9, we present the task category, resource link, encoder module, predictor module, and the employed technical scheme for each method. Unusual abbreviations are explained in the footnote \*.

Sub Category	Method	Encoder	Predictor	Scheme
Functional Property Prediction	LMDisorder [217]	ProtT5	Transformer-based Model	LM-as-Encoder (Repr.)
	ESMtherm [218]	ESM-2	Linear Head	LM-as-Encoder (FM FT)
	VISH-Pred [219]	ESM-2	Ensemble Classifier	LM-as-Encoder (FM FT)
	PIC [220]	ESM-2	Customized Attention Model	LM-as-Encoder (FM FT)
	PeptideBERT [221]	ProtBERT	MLP Head	LM-as-Encoder (FM FT)
	LassoESM [222]	ESM-2	Customized Attention Model	LM-as-Encoder (FM FT)
	OphPred [223]	ESM-2	XGBoost / KNN Classifier	LM-as-Encoder (Repr.)
	EpHod [224]	ESM-1v	Customized Attention Model	LM-as-Encoder (Repr.)
	ProteinGym [160]	ESM-1b / ESM-2 / ProtGPT2 / ESM-MSA-1b / SaProt / ...	-	LM-as-Encoder (Infer.)
	Brandes et al. [225]	ESM-1b	-	LM-as-Encoder (Infer.)
Functional Class Annotation	FSFP [226]	ESM-1v / ESM-2 / SaProt	-	LM-as-Encoder (PE FT & Infer.)
	DePLM [227]	ESM-1v / ESM-2	Denosing Module	LM-as-Encoder (PE FT & Infer.)
	Flamholz et al. [228]	ProtBERT	MLP Head	LM-as-Encoder (Repr.)
	Fernandez et al. [64]	ESM-1b	KNN / SVM / CNN Classifier	LM-as-Encoder (Repr.)
	GraphEC [229]	ProtT5	GNN-based Model & Label Diffusion Algorithm	LM-as-Encoder (Repr.)
	NetGO 3.0 [230]	ESM-1b	Logistic Regression Head	LM-as-Encoder (Repr.)
	DeepGO-SE [231]	ESM-2	Approximate Models	LM-as-Encoder (Repr.)
	SPROF-GO [232]	ProtT5	Label Diffusion Algorithm	LM-as-Encoder (Repr.)
	GraphBepi [233]	ESM-2	GNN & Bi-LSTM-based Model	LM-as-Encoder (Repr.)
	GPSite [66]	ProtT5	GNN-based Model	LM-as-Encoder (Repr.)
Functional Sites Identification	SignalP 6.0 [234]	ProtBERT	CRF Probabilistic Model	LM-as-Encoder (FM FT)
	USPNet [235]	ESM-MSA-1b / ESM-2	Bi-LSTM-based Model	LM-as-Encoder (Repr.)
	ConPlex [67]	ProtBERT & Morgan Fingerprint Encoder	Cosine Distance Calculation	LM-as-Encoder (Repr.)
	BIND [236]	ESM-2 & Ligand Encoder	GNN-based Model	LM-as-Encoder (Repr.)
	UniKP [68]	ProtT5 & SMILES Transformer [237]	Extra Trees Model	LM-as-Encoder (Repr.)
	BALM [238]	ESM-2 & ChemBERTa-2 [239]	MLP-based Model	LM-as-Encoder (PE FT)
	ReactZyme [240]	ESM-2 & Molecule Encoder	GNN & MLP-based Model	LM-as-Encoder (Repr.)
	EasIFA [241]	Enzyme Repr. Branch (w/ ESM-2 & GearNet [242]) & Reaction Repr. Branch	Cross-Attention Module	LM-as-Encoder (Repr.)
	SaLT&PepPr [243]	ESM-2	MLP Head	LM-as-Encoder (PE FT)
	Sledzieski et al. [244]	ESM-2	MLP Head	LM-as-Encoder (PE FT)
Protein-Biomolecule Interactions Prediction	ProLLM [245]	ProtT5	Flan-T5-large [14]	LM-as-Predictor
	Prot2Token [246]	ESM-2 & BARTSmiles [247]	LLM Decoder	LM-as-Predictor
	ProteinChat [33]	GVP-GNN [51] & Projector	Vicuna-13B [248]	LM-as-Predictor
	Prot2Text [70]	ESM-2, RGCN Encoder & Fusion Blocks	GPT-2 [16]	LM-as-Predictor
	ProteinGPT [249]	ESM-2 & GVP-GNN	LLaMA-3 [195]	LM-as-Predictor
	FAPM [250]	ESM-2 & Q-Former Module [89]	Mistral-7B [251]	LM-as-Predictor
	ProtT3 [252]	ESM-2 & Q-Former Module	Galactica [208]	LM-as-Predictor
Multi-Task Question-Answering				

\* **Repr.** - utilization of pLM-produced Representations; **FM FT** - Mull-Model Fine-Tuning; **PE FT** - Parameter-Efficient Fine-Tuning; **Infer.** - Likelihood Inference

substrate structures. While taking AA sequences to represent enzymes and molecular graphs of substrates and products to describe catalyzed reactions, ReactZyme [240] has benchmarked the enzyme-reaction prediction, ranking enzymes by their catalytic ability for specific reactions. In addition, EasIFA [241] is designed to annotate enzyme active sites specific to reactions, predicting the enzymatic activity of residues by aligning the input information of proteins and enzymatic reactions.

Besides, parameter-efficient fine-tuning and prompting techniques have been extensively explored in protein-protein interaction (PPI) prediction. For example, Brixi et al. [243] fine-tune the final few layers of ESM-2 together with a prediction head to identify PPI sites. Sledzieski et al. [244] fine-tune ESM-2 with LoRA for binary PPI prediction, achieving state-of-the-art performance while greatly reducing memory cost. Recently, Protein Chain of Thought (ProCoT) [245] has tried to enhance the reasoning capability of LLM tailored for PPI

prediction, specifically connecting the concepts of signaling pathways and natural language prompts, discovering the interaction between upstream and downstream proteins undergoing multiple biological signal transductions.

#### 4.2.5 Multi-Task Question-Answering

Despite the fantastic performance achieved by the above-mentioned methods, adapting pLMs to each task over and over again can still be time-consuming and computationally expensive. Especially when fresh tasks come up, the consumption of resources will continue to increase in new training rounds. To avoid such trouble, efforts are made to create **multi-task question-answering** frameworks capable of making varied predictions within a unified model, where the *LM-as-Predictor* scheme is typically employed. For example, Prot2Token [246] possess a customized "bi-steam encoders - unified decoder" architecture, integrating a pLM (i.e., ESM-2) and a molecular LM with an autoregressive LLM to perform a variety of protein function predictions. In

addition to the encoded protein and molecule sequences, a specific "task token" is also introduced to the LLM decoder to prompt what predictions should be made, which covers over ten tasks like stability prediction, fluorescence prediction, GO function prediction, human PPI prediction, etc.

As an intuitive and inclusive data form, natural language carries the question-answering process. Given a question and related contexts as prompts, LLMs can generate textual descriptions that encompass broader knowledge rather than fixed labels. Consequently, ChatGPT-like systems are developed to enable users to upload proteins and engage in question-answering interactions to gain insights. For example, Prot2Text [70] combines GNNs and ESM-2 to encode a protein into a fused representation and employs GPT-2 to generate the protein's text description. Protein-Chat [33] consists of a structure encoder (i.e., GVP-GNN), a projection layer, and a general LLM (i.e., Vicuna-13B [248]), which are responsible for producing protein embedding, adapting protein to natural language, and decoding answers, respectively. Similarly, ProteinGPT [249] integrates protein sequence and structure encoders (i.e., ESM-2 and GVP-GNN) with linear projectors for seamless representation adaptation, connected to a general LLM (i.e., LLaMA-3) to generate logically consistent responses. Moreover, both FAPM [250] and ProtT3 [252] employ the Querying Transformer (Q-Former) [89] module to bridge the protein-text modalities, thereby connecting pLM encoder with LLM decoder and achieving accurate protein-to-text generations.

### 4.3 Protein Design

Over millions of years of evolution, natural proteins possess a wide range of functions, which support the running of living systems. However, existing proteins occupy only a minimal fraction of the protein space (illustrated in Figure 1), suggesting that enhanced or even novel functions might be encoded within the extensive unseen protein sequences. Rather than passively watch the slow process of natural evolution unfold, scientists investigate protein design to produce new proteins with desired functions. The critical challenge lies in efficiently exploring the vast protein space to find a manageable number of protein sequences that are plausible, functionally significant, and diverse. Depending on whether starting from existing proteins or scratch, we can classify protein design into two main categories [71]: redesign and *de novo* design.

#### 4.3.1 Protein Redesign

Protein redesign initializes the exploration of protein space from existing proteins, aiming to enhance existing functional properties. Directed evolution [257, 258] is a classic experimental method that emulates the working mechanism of natural evolution in labs, involving an iterative process that alternates between diversification and screening. In each cycle, candidate proteins experience random mutations, followed by experimental measurement of the target functional property (generally the fitness). The most favorable variants are retained as candidates, and the iterative process continues until the desired design goal is achieved [34]. As shown in Figure 10-1, it is considered a sequence optimization process within the fitness landscape, aiming to step from a given seed

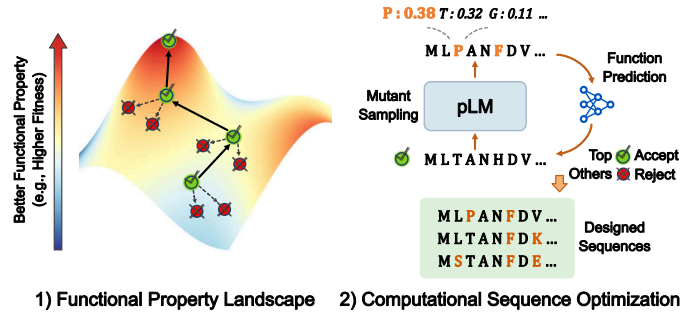


Fig. 10: **Protein Redesign: Function-Oriented Protein Sequence Optimization.** 1) Directed evolution is performed on the functional property landscape over sets of protein mutants. While the landscape exists conceptually, it is generally not fully revealed. Therefore, protein sequence optimization includes iterative rounds of mutant sampling and functional property validation. In each round, only advantaged mutants remain candidates for the next round. 2) Computational protein sequence optimization methods employ pLMs to guide mutant sampling, and leverage the predicted functional property to filter candidate sequences.

(existing sequence) to the optima (functionally enhanced sequence). Nevertheless, the effectiveness and efficiency of this optimization process are limited in two ways: random navigation leads to lots of unnecessary mutants, and repeated measurements of functional properties impose a considerable experimental burden. Consequently, efforts are dedicated to computational sequence optimization to infer superior protein variants, aiming to lessen the laboratory workload and improve the success rate.

In recent years, pLMs have played crucial roles in computational protein redesign, aiding in simulating fitness landscapes and identifying advantageous mutations. Concretely, autoencoding pLMs have learned inter-residue relationships from large-scale pre-training. Taking input as a single wild-type protein sequence  $\mathbf{X} = [x_1, x_2, \dots, x_L]$ , pLMs are expert in producing conditional likelihoods, such as the 'wild-type marginal'  $p(x'_i | \mathbf{X})$  conditioned on the original sequence, or the 'masked marginal'  $p(x'_i | \mathbf{X}_{-i})$  conditioned on the sequence with interested site masked. Those inferred likelihoods should probably adhere to the law of natural protein sequences. Therefore, an important assumption emerges in protein sequence optimization: *mutations with high pLM-inferred likelihoods tend to be evolutionarily plausible* [259].

The distribution of protein sequences learned by pLMs is primarily exploited to assist mutant sampling. Hie et al. [259] perform affinity maturation of human antibodies by employing pLMs to recommend plausible AA substitutions. Ensembled ESM-1b and ESM-1v models compute the 'wild-type marginal' likelihoods of all possible single-residue substitutions on the antibody variable regions, and the substitutions exceeding wild-type likelihood with a consensus are accepted for further analysis. This pLM-guided redesign strategy leads to impressive results, enhancing the affinity of seven antibodies by evaluating only twenty or fewer new variants of each antibody over just two rounds of laboratory evolution. Similarly, Johnson et al. [260] use Gibbs sampling from autoencoding pLMs (e.g., ESM-1b, ProtBERT) to generate novel protein sequences that retain critical characteristics of relevant natural sequences.



The same idea holds for multiple-sequence-based pLMs, which produce conditional likelihoods based on aligned homologous sequences. Sgarbossa et al. [261] propose an iterative protein design method that directly leverages the MLM objective to generate new sequences with ESM-MSA-1b. The resulting sequences are scored at the same level as natural sequences across the homology, co-evolution, and structure-based measures.

In addition, the pLM-guided mutant sampling and function prediction are usually performed collaboratively to contract the range of candidate proteins across iterations. The demonstration can be found in Figure 10-2. Tran et al. [262] utilize pLMs to pinpoint the mutation hotspots and suggest substitutions, then train a fitness prediction model to select top-performance variants. LLM-GA [263] combines pLM with genetic algorithms to optimize enzyme feasibility and turnover dynamics: ESM-2 directs the creation of a pool of mutants, which are screened using specific fitness measurements. EvoOpt [264] leverage ESM-MSA-1b to produce multiple mutant sequences by inferring the randomly masked positions in the original MSA. The generated sequences are subjected to zero-shot fitness prediction, and then the top-ranked proteins are preserved as the proposed candidates. Pro-PRIME [265] is a temperature-guided pLM designed to suggest protein mutants with enhanced stability and activity. It is first utilized for zero-shot mutation suggestion and subsequently fine-tuned for fitness prediction. The eventual top-K mutants are adopted for further experimental analysis. Notably, Darmawan et al. [266] develop a *in silico* evaluation framework that holistically compares pLM-based iterative protein redesign methods. Optimized protein sequences are evaluated for three core criteria: relevance, quality, and diversity, ensuring the designed proteins are functionally relevant and significant while differing from known proteins.

The fundamental goal of protein redesign is to improve the desired functional properties of existing proteins, and the aforementioned sequence optimization strategies achieve this goal by excluding low-fitness mutants. In each iteration, they sample the distribution learned by sequence-based pLMs for evolutionary plausible proteins, which are not controllably biased toward specific functional properties of interest. In enhancing the effectiveness and efficiency of protein redesign, efforts are made to ensure that the mutant sampling process is attuned to desired functions. Conceptionally, function-enhanced pLMs play a role here. As shown in Figure 6-3, Regression Transformer [165] and ProteinNPT [34] learn the bidirectional correlation between protein sequences and functional labels. At a specific masked position, they infer the likelihood over the amino acid vocabulary conditioned on other unmasked tokens and a given functional label. On this basis, we can steer the generation of optimized protein sequences with the property of interest. Moreover, EVOLVEpro [267] is a few-shot active learning framework designed to improve protein activity rapidly, and its effectiveness has been demonstrated by comprehensive experiments across six therapeutically relevant proteins. *In silico*, EVOLVEpro equips ESM-2 with a regression model to learn the activity landscape for each specific protein, thereby guiding the directed evolution process. In each round of directed evolution, a limited number of EVOLVEpro-suggested mutants are evaluated via experimental assays. The data obtained from

these experiments are subsequently used to train the model incrementally and predict mutation candidates for the next round. Over multiple rounds, EVOLVEpro can effectively lead scientists to new proteins with significantly improved desired properties.

Besides, text-guided protein editing emerges as a new computational protein redesign task. Given an initial protein sequence and a prompt describing the desired functional property, we expect to get edited protein sequences that possess the potential to reach our design goal. ProteinDT [166] enables text-guided protein editing based on ProteinCLAP (illustrated in Figure 6-4), which aligns the representation space of protein sequence and textual description. The input protein sequence and text prompt undergo individual encoding and are fused into a latent code, which is further decoded as the optimized protein sequence. Similarly, ProTET [168] comprises two stages: contrastive learning aligns the protein and biotext language model encoders, then the fused representations from original protein sequences and textual instructions serve as the condition for generating edited protein sequences. Moreover, ChatDrug [268] employs conversational LLMs like ChatGPT to edit drugs (e.g., small molecules, peptides, and proteins) through textual descriptions. Users are allowed to ask LLMs to update the drug editing results iteratively. In conclusion, these text-guided protein editing methods demonstrate positive results in computational evaluations, although still away from laboratory validation.

#### 4.3.2 De Novo Protein Design

Rather than engineering existing proteins, *de novo* protein design aims to propose new functional proteins without reference sequences. It is highly challenging as it requires the model to grasp precisely what sequence and structure will achieve the desired function within the vast protein space. Meanwhile, there are distinct advantages of revealing functions never seen in nature and offering complete control over the design progress.

Generally, *de novo* protein design is implemented by reversing the sequence-structure-function paradigm: specify a desired function, design a structure executing this function, and find a sequence that folds into this structure [269]. It is widely acknowledged as two steps: **protein structure generation** and fixed-backbone sequence design (i.e., **inverse folding**) [73]. That is to say, the more conserved protein structures are leveraged to impose greater constraints on exploring the vast sequence space.

- In terms of **protein structure generation**, diffusion models have recently achieved remarkable success [103]. For example, RFDiffusion [74] has demonstrated exceptional performance in generating backbone structures for topology-constrained and unconditional protein design. Chroma [75] is a diffusion-based generative model tailored for proteins and protein complexes. Its generation can be guided toward diverse properties through conditioning settings.
- Regarding **inverse folding**, in addition to graph-based models like ProteinMPNN [77], structure-enhanced pLMs have played a significant role. They can infer reasonable protein sequences while correctly understanding the input structures. As illustrated in Figure 5-2, LM-Design [148] implants a structural adapter into pLM to endow it

with structure awareness. As it is trained to reconstruct the corrupted protein sequence based on the provided backbone structure, LM-Design enables the refinement of protein sequences. During inference, a given structure  $S$  is projected to an initial sequence  $X^{(0)}$ , and then the generated sequence is sampled following Markov process  $X^{(t)} \sim p(X^{(t)}|X^{(t-1)}, S)$  in an iterative manner until some fixed number of steps. Subsequently, ProseLM [149] is another structure-conditioned protein design method based on the adaption of pLMs, where the structural context is incorporated through a set of parameter-efficient adapters. By taking ProGen2, an autoregressive pLM, as the foundation model, proseLM can generate protein sequences autoregressively for a given backbone structure. Similarly, InstructPLM [270] successfully aligns a frozen backbone encoder with a frozen pLM decoder by training a protein structure-sequence adapter. Besides, ESM-IF1 [76] employs GVP layers to extract geometric features, followed by a generic Transformer. ESM-IF1 undergoes supervised training of structure-conditioned autoregressive sequence generation, where the training data is augmented by comprising AlphaFold2 predicted structures.

These two steps could even be combined into an "end-to-end" model. For example, Pinal [271] is a protein design framework that aims to generate protein sequences under the guidance of natural language, with text-to-structure and structure-to-sequence stages seamlessly integrated. The textual description of functions is first translated into structural information through an encoder-decoder network, and then protein sequences are generated based on the structure & description through a re-trained structure-aware pLM.

In recent years, the rise of autoregressive pLMs has introduced new ideas to protein modeling and design. It's feasible to implement **protein sequence generation** directly in unconditional or function-conditioned manners. As the composition of training data can significantly influence the generation distribution of AI models, ensuring alignment between the proteins that models learn and their intended applications is acknowledged as essential for improved protein sequence generation [39]. For example, ProGen2 [105] is primarily designed to generate diverse sequences when pre-trained on universal proteins. Subsequently, fine-tuning enables ProGen2 models for family-specific protein sequence generation. Besides, as shown in Figure 6-2, ProGen [35] and ZymCTRL [272] can propose proteins in response to specific prompt tokens. In addition to zero-shot inference, fine-tuning ProGen using curated tag-sequence instances can enhance its generation of proteins from families with sufficient homologous samples, and fine-tuning ZymCTRL enables the generation of designated proteins that have a higher probability of meeting *in silico* filters and displaying activity similar to their natural counterparts.

Protein design is not completed in a one-way prediction but a complex problem encompassing multiple essential links. In practice, experts would reason with domain knowledge and organize dynamic collaborations between a series of prediction or experimental tools. In addition to working as an independent model, general LLMs like GPT-4 have driven the development of multi-agent systems, bringing a new perspective in promoting protein design. ProtAgent [273] is

an LLM-based platform proposed for *de novo* protein design, where multiple AI agents cooperate within a dynamic environment while each possesses the capability of knowledge retrieval, protein modeling, physics simulation, or results analysis, respectively. ProtAgent demonstrates the potential to minimize the requirement for human interference during the iterative problem-solving process.

## 5 APPLICATIONS

In the field of protein science, the impact of Protein Language Models (pLMs) is exhibited in not only *in silico* modeling but also wet-lab recognized applications. In this section, we discuss some significant application topics, i.e., antibody design, enzyme design, and drug discovery. In each application topic, we concisely show the problem's background, summarize the role of LLMs within a general workflow, and present certain studies with solid experimental analysis.

### 5.1 Antibody Design

Antibody is a type of protein that exists in the immune system, aiding in body defense by identifying and neutralizing harmful entities like bacteria and viruses, also known as antigens [274]. As illustrated in 11-A, central to the recognition process is the interaction between antigen and the complementarity-determining regions (CDRs) of antibody. Then, with antigens neutralized, antibodies successfully help maintain the normal functions of living organisms. Once a life body has defeated an antigen, the antibodies synthesized would remain in the bloodstream, offering protection if the same antigen appears again in the future.

Figure 11-B shows the overall workflow of traditional antibody production. The target antigen is typically introduced into an animal, such as a rabbit or a mouse, and their natural immune response can generate antibodies in the blood. Then, the desired antibodies are obtained with the plasma collected and purified [275]. Despite the massive practical success, antibody production like this is restricted to the innate immunity of animals. On the one hand, it's difficult to control the quality of produced antibodies. On the other hand, wet lab experiments are usually cost-extensively and time-consuming. To overcome these problems, massive efforts are dedicated to artificial antibody design, where pLMs have played an important role in recent years.

As shown in Figure 11-C, pLM-based models can assist antibody design by proposing antibody sequences and structures that specifically bind to the target antigen. These designed antibodies are then tested in computational docking and experimental expression analyses. Their biological effectiveness, especially the antigen-binding affinity, is verified by the expression results in yeast cells. For example, PALM-H3 [276] is proposed for the *de novo* generation of antibody heavy chain CDR 3 (CDR-H3) that meets the desired binding specificity. In evaluations, PALM-H3 can generate antibodies that not only target the SARS-CoV-2 wild-type but also adapt to its emerging variants, including Alpha, Delta, and XBB variants. Besides, Shanker et al. [277] demonstrate that a structure-informed pLM, i.e., ESM-IF can be used to guide the evolution of antibodies. Thirty variants of two therapeutic clinical antibodies are screened for their effectiveness in

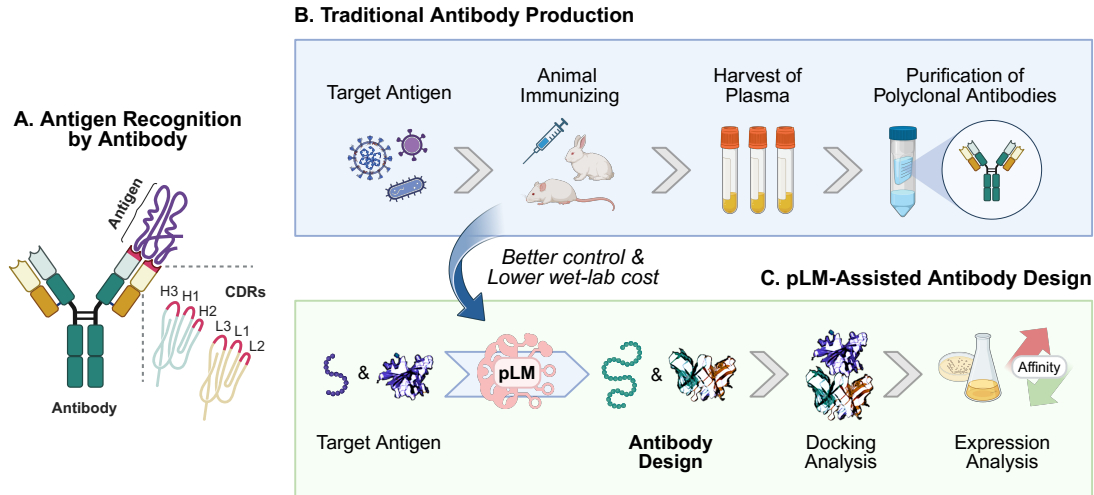


Fig. 11: Overview of **Antibody Design**. (A) In the process of recognizing antigens, the complementarity-determining regions (CDRs) of antibodies are critically important. (B) Traditional antibody production is limited by the immunity of animals. (C) Protein language models can be adopted to propose antibody sequences and structures. The visualization of antibody and antigen structures is derived from PDB entry 7T72, which features SARS-CoV-2 neutralizing antibodies. This figure is created in BioRender.com.

treating severe acute respiratory SARS-CoV-2 infection. In response to antibody-escaped viral variants of concern BQ.1.1 and XBB.1.5, the designed antibodies have demonstrated notable performance, showing up to a 25-fold enhancement in neutralization and a 37-fold improvement in affinity.

## 5.2 Enzyme Design

Enzymes are valuable proteins that act as natural molecular optimization machines to speed up chemical reactions. As illustrated in Figure 12-A, in a suitable environment, enzymes interact with substrates at their active sites, enabling a broad range of biological activities [278]. For example, enzymes break down large molecules like glucose into smaller ones that the body can use for energy [279] and assist in DNA replication [280]. Figure 12-B shows the process of natural enzyme synthesis. Wild-type enzymes are ultimately determined by the information recorded in genes and are produced through the transcription, translation, and folding flows by life bodies. However, wild-type enzymes usually have short lifespans and low stability [281], making it difficult to meet the emerging modern requirements. Therefore, scientists are dedicated to designing enzymes with enhanced properties (e.g., thermostability) and even new catalytic functions, which could help in tackling major global challenges, such as the shortage of energy, pollution of the environment, and lack of sufficient food [282].

Figure 12-C shows a typical workflow of pLM-assisted enzyme design. Scientists can utilize pLMs to optimize wild-type enzymes for desired functions, and the designed enzymes possess the potential to exhibit significant efficiency throughout computational selection and further expression analysis. For example, InstructPLM [270] employs pLMs to generate optimized enzyme sequences based on specific backbone structures. In wet lab experiments (the expression analysis in *E. coli*), it is demonstrated that the redesigned enzymes of PETase and L-MDH exhibit efficacy levels that exceed those of their wild-type. Besides, Johnson et al. [283] conduct computational scoring and experimental assessment

of enzymes generated by three contrasting models. With the help of a novel computational framework named COMPASS, part of the generated sequences with high scores are selected for the further *in vitro* assays. Experimental results demonstrate that certain enzyme sequences designed by pLMs could express in *E. coli* and show activity.

## 5.3 Drug Discovery

In cells, biological functions are usually not performed by individual proteins but rather by their dynamic interactions with each other or with molecules, and the temporary complexes are specifically formed. These interactions are essential for managing the cellular activities essential for life, such as transcription, splicing, etc [284, 285]. Viewed from the pathogenesis perspective, the interactions between drugs (mostly small molecules) and their targets (generally proteins) can effectively regulate the health of organisms. Understanding the target protein and revealing drug-target interaction is critical in drug discovery [286], which could raise great medical significance and economic value.

Figure 13-B outlines a general process of drug discovery. Initially, often in academic settings, scientists would hypothesize that activating or inhibiting a target protein could have therapeutic effects in a disease state. The next step is to confirm if the target exactly plays a crucial role in the onset and progression of diseases. Subsequently, compound screening takes place, aiming to select drugs that can interact with the specific target from a vast library of drug compounds. This process typically involves high-throughput experiments, which require substantial time and materials [287]. After the preliminary screening, secondary analyses are conducted to assess the selected compounds in various biological models and experiments, ensuring their efficacy and safety. Finally, these compounds are tested in animal models to evaluate their pharmacological activity within a biological system.

For decades, researchers have attempted to accelerate drug discovery by computational strategies. Notably, pLM-

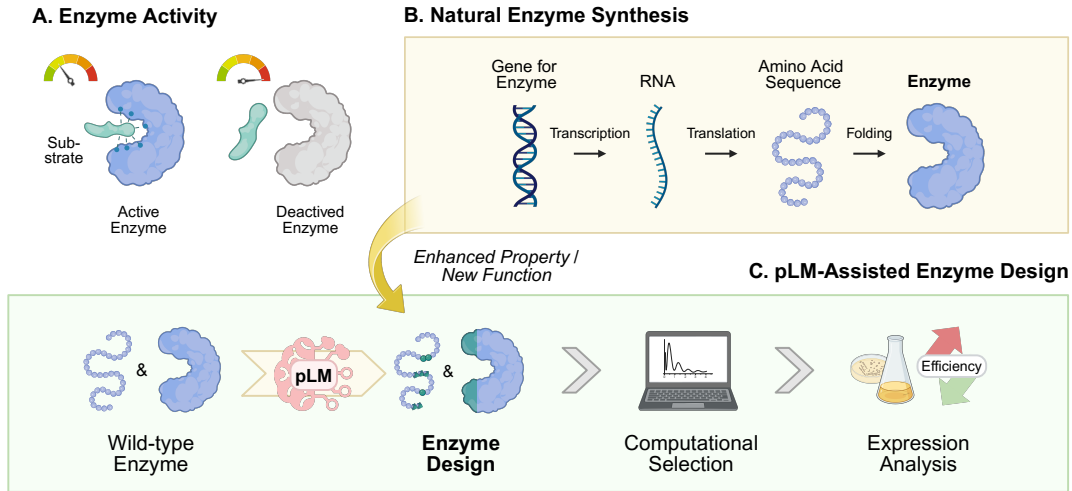


Fig. 12: Overview of **Enzyme Design**. (A) Natural enzymes are usually sensitive to the reaction environment. (B) Genes determine the function of natural enzymes. (C) pLMs can assist in the design of enhanced enzymes. This figure is created in BioRender.com.

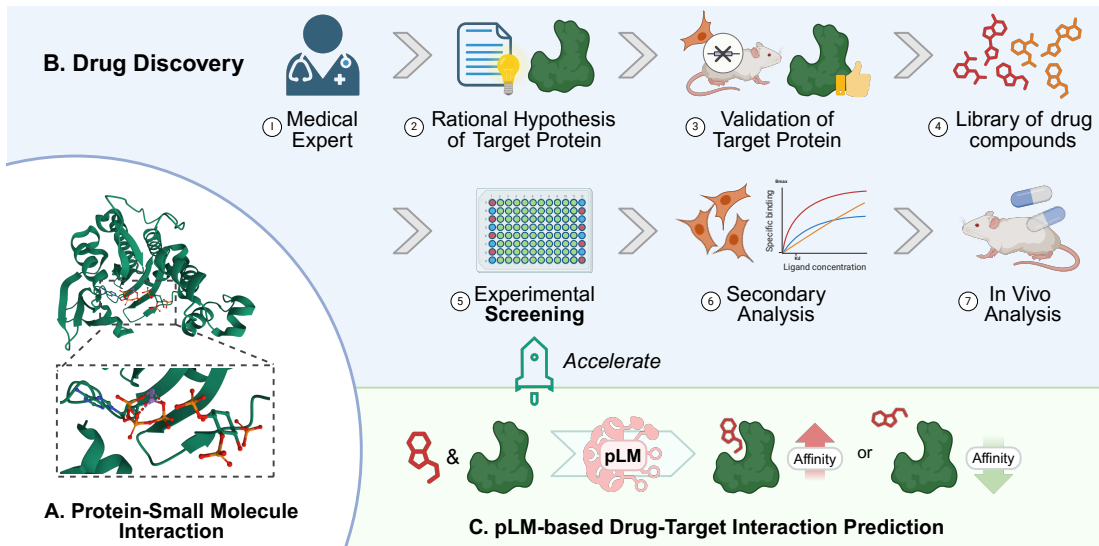


Fig. 13: Overview of **Drug Discovery**. (A) A visualization for protein-ligand interaction, which is obtained from PDB entry 8PPB. (B) Overview of a general drug discovery process [287]. (C) pLMs are employed to predict drug-target interactions, potentially accelerating drug discovery by aiding in drug candidate screening. This figure is created in BioRender.com.

based models demonstrate the potential to complement screening by predicting the interactions between drugs and target proteins. For example, TransDTI [288] is an effective computational workflow that employs pLMs to categorize drug-protein target interactions as active, inactive, or intermediate. With prediction results confirmed through molecular docking and simulation analysis, TransDTI can accurately identify different drug interactions with two proteins, i.e., MAP2k and TGF $\beta$ . In addition, ConPlex [67] investigates contrastive learning in drug-target latent space. The usability of ConPlex is demonstrated by testing 19 kinase-drug interactions, with 12 interactions confirmed. Specifically, four of these interactions demonstrated exceptionally high binding affinity, including a potent EPHB1 inhibitor with a dissociation constant of 1.3 nM when interacting with the compound PD-166326.

## 6 FUTURE DIRECTIONS

In this survey, we have thoroughly reviewed the progress in the interdisciplinary field of computational protein science and large language models, illustrating that protein language models possess a foundational understanding of proteins and can promote the advancement in essential protein modeling problems, i.e., structure prediction, function prediction, and protein design. Despite the achievements, there are still several challenges in this field, which also indicate the future directions.

### 6.1 Data Scarcity

Breakthrough AI for protein studies are primarily driven by abundant protein sequence and structure data. Gratefully, excellent scientists have invested countless efforts in laboratory protein sequencing and structure solutions and have generously shared their findings in public databases.



Moreover, AlphaFold has predicted extensive protein structures that are highly reliable. However, data scarcity is still a tough issue in many specific practical tasks. For example, the imbalanced species representations in protein sequence databases lead to consistent species bias that can be detrimental for protein design applications [253]; and it's difficult to assess the robustness of models developed for protein fitness prediction and redesign as lacking large-scale benchmarks with consistent ground-truth labels [160]. The enhancement of multi-modal learning is also limited by the scale of protein-text datasets developed with expertise [88]. To address the challenge of data scarcity, it's promising to expand and diversify the training data through augmentation or synthesis methods [190], and empower AI models with the capability to learn from limited instances.

## 6.2 Protein Interaction Modeling

In living organisms, proteins usually do not exist in a simple form of "single-chain sequence, monomer structure, & independent function". Instead, multiple polypeptide chains could fold together to form protein complexes, and various protein molecules can interact physically within specific biomolecular contexts. These protein-protein interactions are crucial for regulating a wide array of biological activities. Despite the current success of language model techniques in computational protein science, extending protein language models (pLMs) to effectively model protein interactions remains a significant challenge, which involves foundational tasks such as understanding, structure prediction, and function prediction of protein complexes. One of the primary reasons is the limited availability of data on protein interactions. For instance, the Protein Data Bank (PDB) [45] contains relatively few protein complexes, and the Observed Antibody Space (OAS) [134] has limited paired variable heavy and light chain data. Meanwhile, most existing protein interaction data may not have been fully exploited in training pLMs. For instance, data pipelines, such as those employed by AlphaFold2 [8] and ESMFold [25], often treat polypeptide chains individually, even when they are components of larger protein complexes. Navigating the challenge alongside the opportunity, several efforts have been made to advance protein interaction modeling. For example, IgBert and IgT5 [289] are antibody-specific language models designed to consistently handle the paired and unpaired variable region sequences. Linker-Tuning [290] extends ESMFold, a method originally designed for single-chain structure prediction, to predict heterodimer structures. Furthermore, RDE [291] focuses on understanding the effect of mutations on protein-protein interactions. In the future, we believe that there will be a surge in research and novel findings aimed at addressing this formidable challenge. Continued efforts in this field hold great promise for enhancing our understanding of complex biological systems and advancing computational protein science.

## 6.3 Explainability

As summarized in Sections 3 and 4, pLMs extract statistical patterns within natural protein sequences, and then pLM-based structure prediction, function prediction, and protein design methods capture information flows within

the sequence-structure-function paradigm. However, the learned principles are latent in hundreds of millions of model parameters. Although AI models such as ESM-2 and AlphaFold2 have been widely acknowledged and applied in scientific exploration, scientists have barely gained explicit physical insights from them [292]. Experts would take "black-box" prediction models as tools that output probably accurate "observations", while do not directly get aids in the understanding of new theories. If we can distill the knowledge learned by AI models in a form humans can readily understand, protein sciences could be significantly promoted in a new manner. These days, there have been individual studies on this. Zhang et al. [293] demonstrate that pLMs learn evolutionary statistics of interacting sequence motifs. Then, InterPro [294] is proposed to extract, analyze, and visualize human-interpretable latent features from pLMs. In the future, the more in-depth Explainable AI (XAI) [295] research for protein science still presents a challenging but influential blue ocean.

## 6.4 Bridging Computational and Experimental Research

Scientific discovery is a multifaceted process that involves multiple interconnected stages covering hypothesis, experiments, and data [296]. To date, protein language models are mainly accelerating protein science research with dry lab predictions. Biologists still undertake heavy workloads like rational reasoning with domain knowledge, planning the combination of computational tools & experimental facilities, and performing wet lab experiments. It is acknowledged that general LLMs like GPT-4 [297] have exceptional abilities in solving complex problems by empowering multiple intelligence agents, and that could be utilized in driving end-to-end scientific research. Then, there have been some surprising results in the field of chemistry. Artificial intelligence systems empowered by GPT-4, like ChemCrow [298] and Coscientist [299], can autonomously design, plan, and perform complex experiments by incorporating multiple scientific research tools. In a similar way, LLMs should hold the potential to drive further research in protein science [273]. In the future, we look forward to more mature studies in autonomous research that bridge the gap between computational and experimental protein science, thereby helping biologists in all stages of work and accelerating real-world scientific discovery.

## 6.5 Computational Efficiency

In the recent progress of AI, scaling law [81] stands out as a particularly instructive concept. As LLMs scale up in parameters, data, and computing, larger models acquire unprecedented emergent capabilities and demonstrate substantially improved performance. It is observed the development of pLMs follows this trend as well. For example, ESM-3 [30] is scaled up to 98 billion parameters, and xTrimoPGLM is trained at a massive scale of 100 billion parameters. No matter for big companies or academic groups, the high computational cost poses a tough issue, and the computational efficiency of pLMs should be improved. Therefore, it's in demand to understand the scaling behaviors for protein language modeling and formulate optimized computational schemes. Recently, there have been certain

ongoing studies questioning whether pLMs are compute-optimal and reached the unanimous conclusion of non-optimal [300, 301], exploring how to scale down pLMs for better efficiency while maintaining their expressiveness [302, 303], and implementing novel techniques like FlashAttention [304] to achieve efficient inference, training, and fine-tuning of pLMs [305, 306]. In addition to these preliminary explorations, in the future, more efforts could be dedicated to implementing optimal models in the constraint of predetermined computation budgets.

## 7 CONCLUSION

We present a comprehensive survey of computational protein science in the era of LLMs, containing broad content from background concepts to the latest advancements. First, we outline the biological basis and data profiles in protein modeling. Second, we review three categories of pLMs with abilities to comprehend amino acid sequences, recognize structural and functional information, and bridge multiple biomedical languages. Next, we introduce the utilization and adaptation of pLMs, highlighting their significant impacts on structure prediction, function prediction, and protein design. Then, we specify the application potentials of pLMs in antibody design, enzyme design, and drug target discovery. Finally, we share the promising future directions in this fast-growing field.

## REFERENCES

- [1] Noelia Ferruz and Birte Höcker. Controllable protein design with language models. *Nature Machine Intelligence*, 4(6):521–532, 2022.
- [2] Tristan Beppler and Bonnie Berger. Learning the protein language: Evolution, structure, and function. *Cell systems*, 12(6):654–669, 2021.
- [3] Christian B Anfinsen. Principles that govern the folding of protein chains. *Science*, 181(4096):223–230, 1973.
- [4] Oliver C Redfern, Benoit Dessailly, and Christine A Orengo. Exploring the structure and function paradigm. *Current opinion in structural biology*, 18(3):394–402, 2008.
- [5] John Maynard Smith. Natural selection and the concept of a protein space. *Nature*, 225(5232):563–564, 1970.
- [6] Dan Ofer, Nadav Brandes, and Michal Linial. The language of proteins: Nlp, machine learning & protein sequences. *Computational and Structural Biotechnology Journal*, 19:1750–1758, 2021.
- [7] Ethan C Alley, Grigory Khimulya, Surojit Biswas, Mohammed AlQuraishi, and George M Church. Unified rational protein engineering with sequence-based deep representation learning. *Nature methods*, 16(12):1315–1322, 2019.
- [8] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [9] Minkyung Baek, Frank DiMaio, Ivan Anishchenko, Justas Daurapas, Sergey Ovchinnikov, Gyu Rie Lee, Jue Wang, Qian Cong, Lisa N Kinch, R Dustin Schaeffer, et al. Accurate prediction of protein structures and interactions using a three-track neural network. *Science*, 373(6557):871–876, 2021.
- [10] Maxat Kulmanov and Robert Hoehndorf. Deepgopplus: improved protein function prediction from sequence. *Bioinformatics*, 36(2):422–429, 2020.
- [11] Lukas Theo Schmitt, Maciej Paszkowski-Rogacz, Florian Jug, and Frank Buchholz. Prediction of designer-recombinases for dna editing with generative deep learning. *Nature Communications*, 13(1):7966, 2022.
- [12] Raphael R Eguchi, Christian A Choe, and Po-Ssu Huang. Ig-vae: Generative modeling of protein structure by direct 3d coordinate generation. *PLoS computational biology*, 18(6):e1010271, 2022.
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [14] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- [15] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [16] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [17] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [18] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [19] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [20] Zihuai Zhao, Wenqi Fan, Jiatong Li, Yunqing Liu, Xiaowei Mei, Yiqi Wang, Zhen Wen, Fei Wang, Xiangyu Zhao, Jiliang Tang, et al. Recommender systems in the era of large language models (llms). *IEEE Transactions on Knowledge and Data Engineering*, 2024.
- [21] Haohao Qu, Wenqi Fan, Zihuai Zhao, and Qing Li. Tokenrec: Learning to tokenize id for llm-based generative recommendation. *arXiv preprint arXiv:2406.10450*, 2024.
- [22] Kai He, Rui Mao, Qika Lin, Yucheng Ruan, Xiang Lan, Mengling Feng, and Erik Cambria. A survey of large language models for healthcare: from data, technology, and applications to accountability and ethics. *arXiv preprint arXiv:2310.05694*, 2023.
- [23] Micaela E Consens, Cameron Dufault, Michael Wainberg, Duncan Forster, Mehran Karimzadeh, Hani Goodarzi, Fabian J Theis, Alan Moses, and Bo Wang. To transformers and beyond: Large language models for the genome. *arXiv preprint arXiv:2311.07621*, 2023.
- [24] Jiatong Li, Yunqing Liu, Wenqi Fan, Xiao-Yong Wei, Hui Liu, Jiliang Tang, and Qing Li. Empowering molecule discovery for molecule-caption translation with large language models: A chatgpt perspective. *IEEE Transactions on Knowledge and Data Engineering*, 2024.
- [25] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- [26] Noelia Ferruz, Steffen Schmidt, and Birte Höcker. Protgpt2 is a deep unsupervised language model for protein design. *Nature communications*, 13(1):4348, 2022.
- [27] Bo Chen, Xingyi Cheng, Pan Li, Yangli-ao Geng, Jing Gong, Shen Li, Zhilei Bei, Xu Tan, Boyan Wang, Xin Zeng, et al. xtrimopglm: unified 100b-scale pre-trained transformer for deciphering the language of protein. *arXiv preprint arXiv:2401.06199*, 2024.
- [28] Roshan M Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu, and Alexander Rives. Msa transformer. In *International Conference on Machine Learning*, pages 8844–8856. PMLR, 2021.
- [29] Jin Su, Chencheng Han, Yuyang Zhou, Junjie Shan, Xibin Zhou, and Fajie Yuan. Saprot: protein language modeling with structure-aware vocabulary. *bioRxiv*, pages 2023–10, 2023.
- [30] Thomas Hayes, Roshan Rao, Halil Akin, Nicholas J Sofroniew, Deniz Oktay, Zeming Lin, Robert Verkuil, Vincent Q Tran, Jonathan Deaton, Marius Wiggert, et al. Simulating 500 million years of evolution with a language model. *Science*, page eads0018, 2025.
- [31] Liuzhenghao Lv, Zongying Lin, Hao Li, Yuyang Liu, Jiayi Cui, Calvin Yu-Chian Chen, Li Yuan, and Yonghong Tian. Prollama: A protein large language model for multi-task protein language processing. *arXiv preprint arXiv:2402.16445*, 2024.
- [32] Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. *arXiv preprint arXiv:2310.07276*, 2023.

- [33] Han Guo, Mingjia Huo, Ruiyi Zhang, and Pengtao Xie. Proteinchat: Towards achieving chatgpt-like functionalities on protein 3d structures. *Authorea Preprints*, 2023.
- [34] Pascal Notin, Ruben Weitzman, Debora Marks, and Yarin Gal. Proteinnp: improving protein property prediction and design with non-parametric transformers. *Advances in Neural Information Processing Systems*, 36, 2024.
- [35] Ali Madani, Ben Krause, Eric R Greene, Subu Subramanian, Benjamin P Mohr, James M Holton, Jose Luis Olmos, Caiming Xiong, Zachary Z Sun, Richard Socher, et al. Large language models generate functional protein sequences across diverse families. *Nature Biotechnology*, 41(8):1099–1106, 2023.
- [36] Qiang Zhang, Keyang Ding, Tianwen Lyv, Xinda Wang, Qingyu Yin, Yiwen Zhang, Jing Yu, Yuhao Wang, Xiaotong Li, Zhuoyi Xiang, et al. Scientific large language models: A survey on biological & chemical domains. *arXiv preprint arXiv:2401.14656*, 2024.
- [37] Bozhen Hu, Jun Xia, Jiangbin Zheng, Cheng Tan, Yufei Huang, Yongjie Xu, and Stan Z Li. Protein language models and structure prediction: Connection and progression. *arXiv preprint arXiv:2211.16742*, 2022.
- [38] Serbulent Unsal, Heval Atas, Muammer Albayrak, Kemal Turhan, Aybar C Acar, and Tunca Doğan. Learning functional properties of proteins with language models. *Nature Machine Intelligence*, 4(3):227–245, 2022.
- [39] Jeffrey A Ruffolo and Ali Madani. Designing proteins with language models. *Nature Biotechnology*, 42(2):200–202, 2024.
- [40] Edmond De Hoffmann and Vincent Stroobant. *Mass spectrometry: principles and applications*. John Wiley & Sons, 2007.
- [41] Randy J Read, Paul D Adams, W Bryan Arendall, Axel T Brunger, Paul Emsley, Robbie P Joosten, Gerard J Kleywegt, Eugene B Krissinel, Thomas Lütke, Zbyszek Otwinowski, et al. A new generation of crystallographic validation tools for the protein data bank. *Structure*, 19(10):1395–1412, 2011.
- [42] Richard Henderson, Andrej Sali, Matthew L Baker, Bridget Carragher, Batsal Devkota, Kenneth H Downing, Edward H Egelman, Zukang Feng, Joachim Frank, Nikolaus Grigorieff, et al. Outcome of the first electron microscopy validation task force meeting. *Structure*, 20(2):205–214, 2012.
- [43] Douglas M Fowler and Stanley Fields. Deep mutational scanning: a new style of protein science. *Nature methods*, 11(8):801–807, 2014.
- [44] Baris E Suzek, Yuqi Wang, Hongzhan Huang, Peter B McGarvey, Cathy H Wu, and UniProt Consortium. Uniref clusters: a comprehensive and scalable alternative for improving sequence similarity searches. *Bioinformatics*, 31(6):926–932, 2015.
- [45] Protein data bank: the single global archive for 3d macromolecular structure data. *Nucleic acids research*, 47(D1):D520–D528, 2019.
- [46] Michael Ashburner, Catherine A Ball, Judith A Blake, David Botstein, Heather Butler, J Michael Cherry, Allan P Davis, Kara Dolinski, Selina S Dwight, Janan T Eppig, et al. Gene ontology: tool for the unification of biology. *Nature genetics*, 25(1):25–29, 2000.
- [47] Francis Crick. Central dogma of molecular biology. *Nature*, 227(5258):561–563, 1970.
- [48] Eugene V Koonin and Artem S Novozhilov. Origin and evolution of the genetic code: the universal enigma. *IUBMB life*, 61(2):99–111, 2009.
- [49] Cyrus Chothia. One thousand families for the molecular biologist. *Nature*, 357(6379), 1992.
- [50] Ken A Dill and Justin L MacCallum. The protein-folding problem, 50 years on. *science*, 338(6110):1042–1046, 2012.
- [51] Bowen Jing, Stephan Eismann, Patricia Suriana, Raphael John Lamarre Townshend, and Ron Dror. Learning from protein structure with geometric vector perceptrons. In *International Conference on Learning Representations*, 2020.
- [52] Zhangyang Gao, Cheng Tan, Pablo Chacón, and Stan Z Li. Pifold: Toward effective and efficient protein inverse folding. *arXiv preprint arXiv:2209.12643*, 2022.
- [53] Philip A Romero and Frances H Arnold. Exploring protein fitness landscapes by directed evolution. *Nature reviews Molecular cell biology*, 10(12):866–876, 2009.
- [54] Lirong Wu, Yufei Huang, Haitao Lin, and Stan Z Li. A survey on protein representation learning: Retrospect and prospect. *arXiv preprint arXiv:2301.00813*, 2022.
- [55] Kevin K Yang, Nicolo Fusi, and Alex X Lu. Convolutions are competitive with transformers for protein sequence pretraining. *Cell Systems*, 15(3):286–294, 2024.
- [56] Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Yu Wang, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, et al. Prottrans: Toward understanding the language of life through self-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 44(10):7112–7127, 2021.
- [57] Zuobai Zhang, Minghao Xu, Arian Jamasb, Vijil Chenthamarakshan, Aurelie Lozano, Payel Das, and Jian Tang. Protein representation learning by geometric structure pretraining. *arXiv preprint arXiv:2203.06125*, 2022.
- [58] Zichen Wang, Steven A Combs, Ryan Brand, Miguel Romero Calvo, Panpan Xu, George Price, Nataliya Golovach, Emmanuel O Salawu, Colby J Wise, Sri Priya Ponnappalli, et al. Lm-gvp: an extensible sequence and structure informed deep learning framework for protein property prediction. *Scientific reports*, 12(1):6832, 2022.
- [59] Z Zhang, C Wang, M Xu, V Chenthamarakshan, AC Lozano, P Das, and J Tang. A systematic study of joint representation learning on protein sequences and structures. *Preprint at http://arxiv.org/abs/2303.06275*, 2023.
- [60] Milot Mirdita, Konstantin Schütze, Yoshitaka Moriawaki, Lim Heo, Sergey Ovchinnikov, and Martin Steinegger. Colabfold: making protein folding accessible to all. *Nature methods*, 19(6):679–682, 2022.
- [61] Wenkai Wang, Zhenling Peng, and Jianyi Yang. Single-sequence protein structure prediction using supervised transformer protein language models. *Nature Computational Science*, 2(12):804–814, 2022.
- [62] Xiaomin Fang, Fan Wang, Lihang Liu, Jingzhou He, Dayong Lin, Yingfei Xiang, Kunrui Zhu, Xiaonan Zhang, Hua Wu, Hui Li, et al. A method for multiple-sequence-alignment-free protein structure prediction using a protein language model. *Nature Machine Intelligence*, 5(10):1087–1096, 2023.
- [63] Christian Dallago, Jody Mou, Kadina E Johnston, Bruce J Wittmann, Nicholas Bhattacharya, Samuel Goldman, Ali Madani, and Kevin K Yang. Flip: Benchmark tasks in fitness landscape inference for proteins. *bioRxiv*, pages 2021–11, 2021.
- [64] Diego Fernández, Álvaro Olivera-Nappa, Roberto Uribe-Paredes, and David Medina-Ortiz. Exploring machine learning algorithms and protein language models strategies to develop enzyme classification systems. In *International Work-Conference on Bioinformatics and Biomedical Engineering*, pages 307–319. Springer, 2023.
- [65] Yufan Liu and Boxue Tian. Protein–dna binding sites prediction based on pre-trained protein language model and contrastive learning. *Briefings in Bioinformatics*, 25(1):bbad488, 2024.
- [66] Qianmu Yuan, Chong Tian, and Yuedong Yang. Genome-scale annotation of protein binding sites via language model and geometric deep learning. *Elife*, 13:RP93695, 2024.
- [67] Rohit Singh, Samuel Sledzieski, Bryan Bryson, Lenore Cowen, and Bonnie Berger. Contrastive learning in protein language space predicts interactions between drugs and protein targets. *Proceedings of the National Academy of Sciences*, 120(24):e2220778120, 2023.
- [68] Han Yu, Huaxiang Deng, Jiahui He, Jay D Keasling, and Xiaozhou Luo. Unikp: a unified framework for the prediction of enzyme kinetic parameters. *Nature Communications*, 14(1):8211, 2023.
- [69] Chao Wang, Hehe Fan, Ruijie Quan, and Yi Yang. Protchatgpt: Towards understanding proteins with large language models. *arXiv preprint arXiv:2402.09649*, 2024.
- [70] Hadi Abdine, Michail Chatzianastasis, Costas Bouyioukos, and Michalis Vazirgiannis. Prot2text: Multimodal protein’s function generation with gnn and transformers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 10757–10765, 2024.
- [71] Pascal Notin, Nathan Rollins, Yarin Gal, Chris Sander, and Debora Marks. Machine learning for functional protein design. *Nature Biotechnology*, 42(2):216–228, 2024.
- [72] Joshua Meier, Roshan Rao, Robert Verkuil, Jason Liu, Tom Sercu, and Alex Rives. Language models enable zero-shot prediction of the effects of mutations on protein function. *Advances in neural information processing systems*, 34:29287–29303, 2021.
- [73] Tanja Kortemme. De novo protein design—from new structures to programmable functions. *Cell*, 187(3):526–544, 2024.
- [74] Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Trippe, Jason Yim, Helen E Eisenach, Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, et al. De novo design of protein structure and function with rdiffusion. *Nature*, 620(7976):1089–1100, 2023.

- [75] John B Ingraham, Max Baranov, Zak Costello, Karl W Barber, Wujie Wang, Ahmed Ismail, Vincent Frappier, Dana M Lord, Christopher Ng-Thow-Hing, Erik R Van Vlack, et al. Illuminating protein space with a programmable generative model. *Nature*, 623(7989):1070–1078, 2023.
- [76] Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer, and Alexander Rives. Learning inverse folding from millions of predicted structures. In *International conference on machine learning*, pages 8946–8970. PMLR, 2022.
- [77] Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J Ragotte, Lukas F Milles, Basile IM Wicky, Alexis Courbet, Rob J de Haas, Neville Bethel, et al. Robust deep learning-based protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022.
- [78] Sarah Alamdari, Nitya Thakkar, Rianne van den Berg, Alex Xijie Lu, Nicolo Fusi, Ava Pardis Amini, and Kevin K Yang. Protein generation with evolutionary diffusion: sequence is all you need. *bioRxiv*, pages 2023–09, 2023.
- [79] Wenqi Fan, Shijie Wang, Jiani Huang, Zhikai Chen, Yu Song, Wenzhuo Tang, Haitao Mao, Hui Liu, Xiaorui Liu, Dawei Yin, et al. Graph machine learning in the era of large language models (llms). *arXiv preprint arXiv:2404.14928*, 2024.
- [80] Yuyan Liu, Sirui Ding, Sheng Zhou, Wenqi Fan, and Qiaoyu Tan. Moleculargpt: Open large language model (llm) for few-shot molecular property prediction. *arXiv preprint arXiv:2406.12950*, 2024.
- [81] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [82] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [83] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.
- [84] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [85] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- [86] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [87] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- [88] Qizhi Pei, Lijun Wu, Kaiyuan Gao, Jinhua Zhu, Yue Wang, Zun Wang, Tao Qin, and Rui Yan. Leveraging biomolecule and natural language through multi-modal learning: A survey. *arXiv preprint arXiv:2403.01528*, 2024.
- [89] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023.
- [90] Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 6491–6501, 2024.
- [91] Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15):e2016239118, 2021.
- [92] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [93] ESM Team. ESM Cambrian: Revealing the mysteries of proteins with unsupervised learning. <https://evolutionaryscale.ai/blog/esm-cambrian>, December 2024.
- [94] Liang He, Shizhuo Zhang, Lijun Wu, Huanhuan Xia, Fusong Ju, He Zhang, Siyuan Liu, Yingce Xia, Jianwei Zhu, Pan Deng, et al. Pre-training co-evolutionary protein representation via a pairwise masked language model. *arXiv preprint arXiv:2110.15527*, 2021.
- [95] Yingheng Wang, Zichen Wang, Gil Sadeh, Luca Zancato, Alessandro Achille, George Karypis, and Huzefa Rangwala. Long-context protein language model. *bioRxiv*, pages 2024–10, 2024.
- [96] Haohao Qu, Liangbo Ning, Rui An, Wenqi Fan, Tyler Derr, Xin Xu, and Qing Li. A survey of mamba. *arXiv preprint arXiv:2408.01129*, 2024.
- [97] Haohao Qu, Yifeng Zhang, Liangbo Ning, Wenqi Fan, and Qing Li. Ssd4rec: a structured state space duality model for efficient sequential recommendation. *arXiv preprint arXiv:2409.01192*, 2024.
- [98] Lei Wang, Hui Zhang, Wei Xu, Zhidong Xue, and Yan Wang. Deciphering the protein landscape with protflash, a lightweight language model. *Cell Reports Physical Science*, 4(10), 2023.
- [99] Yaron Geffen, Yanay Ofran, and Ron Unger. Distilprotbert: a distilled protein language model used to distinguish between real proteins and their randomly shuffled counterparts. *Bioinformatics*, 38(Supplement\_2):ii95–ii98, 2022.
- [100] Ning Sun, Shuxian Zou, Tianhua Tao, Sazan Mahbub, Dian Li, Yonghao Zhuang, Hongyi Wang, Xingyi Cheng, Le Song, and Eric P Xing. Mixture of experts enable efficient and effective protein understanding and design. *bioRxiv*, pages 2024–11, 2024.
- [101] Le Song, Eran Segal, and Eric Xing. Toward ai-driven digital organism: Multiscale foundation models for predicting, simulating and programming biology at all levels. *arXiv preprint arXiv:2412.06993*, 2024.
- [102] Xinyou Wang, Zaixiang Zheng, Fei Ye, Dongyu Xue, Shujian Huang, and Quanquan Gu. Diffusion language models are versatile protein learners. *arXiv preprint arXiv:2402.18567*, 2024.
- [103] Chengyi Liu, Wenqi Fan, Yunqing Liu, Jiatong Li, Hang Li, Hui Liu, Jiliang Tang, and Qing Li. Generative diffusion models on graphs: methods and applications. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, pages 6702–6711, 2023.
- [104] Daniel Hesslow, Niccolò Zanichelli, Pascal Notin, Iacopo Poli, and Debora Marks. Rita: a study on scaling up generative protein sequence models. *arXiv preprint arXiv:2205.05789*, 2022.
- [105] Erik Nijkamp, Jeffrey A Ruffolo, Eli N Weinstein, Nikhil Naik, and Ali Madani. Progen2: exploring the boundaries of protein language models. *Cell systems*, 14(11):968–978, 2023.
- [106] Ahmed Elnaggar, Hazem Essam, Wafaa Salah-Eldin, Walid Moustafa, Mohamed Elkerdawy, Charlotte Rochereau, and Burkhard Rost. Ankh: Optimized protein language model unlocks general-purpose modelling. *arXiv preprint arXiv:2301.06568*, 2023.
- [107] Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. Glm: General language model pretraining with autoregressive blank infilling. *arXiv preprint arXiv:2103.10360*, 2021.
- [108] Jeffrey A Ruffolo, Jeffrey J Gray, and Jeremias Sulam. Deciphering antibody affinity maturation with language models and weakly supervised learning. *arXiv preprint arXiv:2112.07782*, 2021.
- [109] Jinwoo Leem, Laura S Mitchell, James HR Farmery, Justin Barton, and Jacob D Galson. Deciphering the language of antibodies using self-supervised learning. *Patterns*, 3(7), 2022.
- [110] Tobias H Olsen, Iain H Moal, and Charlotte M Deane. Ablang: an antibody language model for completing antibody sequences. *Bioinformatics Advances*, 2(1):vbac046, 2022.
- [111] Kaiyuan Gao, Lijun Wu, Jinhua Zhu, Tianbo Peng, Yingce Xia, Liang He, Shufang Xie, Tao Qin, Haiguang Liu, Kun He, et al. Pre-training antibody language models for antigen-specific computational antibody design. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 506–517, 2023.
- [112] Igor Melnyk, Vijil Chenthamarakshan, Pin-Yu Chen, Payel Das, Amit Dhurandhar, Inkit Padhi, and Devleena Das. Reprogramming pretrained language models for antibody sequence infilling. In *International Conference on Machine Learning*, pages 24398–24419. PMLR, 2023.
- [113] Richard W Shuai, Jeffrey A Ruffolo, and Jeffrey J Gray. Iglm: Infilling language modeling for antibody sequence design. *Cell Systems*, 14(11):979–989, 2023.
- [114] Oliver Marcus Turnbull, Dino Oglic, Rebecca Croasdale-Wood, and Charlotte M Deane. p-iggen: A paired antibody generative language model. *bioRxiv*, pages 2024–08, 2024.



- [115] Simon KS Chu and Kathy Y Wei. Generative antibody design for complementary chain pairing sequences through encoder-decoder language model. *arXiv preprint arXiv:2301.02748*, 2023.
- [116] Michael Remmert, Andreas Biegert, Andreas Hauser, and Johannes Söding. Hhblits: lightning-fast iterative protein sequence searching by hmm-hmm alignment. *Nature methods*, 9(2):173–175, 2012.
- [117] Martin Steinegger and Johannes Söding. Mmseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nature biotechnology*, 35(11):1026–1028, 2017.
- [118] Soumya Ram and Tristan Bepler. Few shot protein generation. *arXiv preprint arXiv:2204.01168*, 2022.
- [119] Le Zhang, Jiayang Chen, Tao Shen, Yu Li, and Siqi Sun. Enhancing the protein tertiary structure prediction by multiple sequence alignment generation. *arXiv preprint arXiv:2306.01824*, 2023.
- [120] Bo Chen, Zhilei Bei, Xingyi Cheng, Pan Li, Jie Tang, and Le Song. Msagpt: Neural prompting protein structure prediction via msa generative pre-training. *arXiv preprint arXiv:2406.05347*, 2024.
- [121] Timothy Truong Jr and Tristan Bepler. Poet: A generative model of protein families as sequences-of-sequences. *Advances in Neural Information Processing Systems*, 36, 2024.
- [122] Damiano Sgarbossa, Cyril Malbranke, and Anne-Florence Bitbol. Protmamba: a homology-aware but alignment-free protein state space model. *bioRxiv*, pages 2024–05, 2024.
- [123] Carlos Outeiral and Charlotte M Deane. Codon language embeddings provide strong signals for use in protein engineering. *Nature Machine Intelligence*, 6(2):170–179, 2024.
- [124] Logan Hallee, Nikolaos Rafailidis, and Jason P Gleghorn. cdsbert-extending protein language models with codon awareness. *bioRxiv*, 2023.
- [125] Zhangzhi Peng, Benjamin Schussheim, and Pranam Chatterjee. Ptm-mamba: A ptm-aware protein language model with bidirectional gated mamba blocks. *bioRxiv*, pages 2024–02, 2024.
- [126] Lorna Richardson, Ben Allen, Germana Baldi, Martin Bera-cochea, Maxwell L Bileschi, Tony Burdett, Josephine Burgin, Juan Caballero-Pérez, Guy Cochrane, Lucy J Colwell, et al. Mgnify: the microbiome sequence data analysis resource in 2023. *Nucleic Acids Research*, 51(D1):D753–D759, 2023.
- [127] I-Min A Chen, Ken Chu, Krishnaveni Palaniappan, Anna Ratner, Jinghua Huang, Marcel Huntemann, Patrick Hajek, Stephan J Ritter, Cody Webb, Dongying Wu, et al. The img/m data management and analysis system v. 7: content updates and new features. *Nucleic acids research*, 51(D1):D723–D732, 2023.
- [128] Martin Steinegger, Milot Mirdita, and Johannes Söding. Protein-level assembly increases protein sequence recovery from metagenomic samples manyfold. *Nature methods*, 16(7):603–606, 2019.
- [129] Z Lan. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.
- [130] K Clark. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555*, 2020.
- [131] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- [132] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- [133] Zhilin Yang. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv preprint arXiv:1906.08237*, 2019.
- [134] Aleksandr Kovaltsuk, Jinwoo Leem, Sebastian Kelm, James Snowden, Charlotte M Deane, and Konrad Krawczyk. Observed antibody space: a resource for data mining next-generation sequencing of antibody repertoires. *The Journal of Immunology*, 201(8):2502–2509, 2018.
- [135] James Dunbar, Konrad Krawczyk, Jinwoo Leem, Terry Baker, Angelika Fuchs, Guy Georges, Jiye Shi, and Charlotte M Deane. Sabdab: the structural antibody database. *Nucleic acids research*, 42(D1):D1140–D1146, 2014.
- [136] Jaina Mistry, Sara Chuguransky, Lowri Williams, Matloob Qureshi, Gustavo A Salazar, Erik LL Sonnhammer, Silvio CE Tosatto, Lisanna Paladin, Shriya Raj, Lorna J Richardson, et al. Pfam: The protein families database in 2021. *Nucleic acids research*, 49(D1):D412–D419, 2021.
- [137] Milot Mirdita, Lars Von Den Driesch, Clovis Galiez, Maria J Martin, Johannes Söding, and Martin Steinegger. Uniclust databases of clustered and deeply annotated protein sequences and alignments. *Nucleic acids research*, 45(D1):D170–D176, 2017.
- [138] Gustaf Ahdriz, Nazim Bouatta, Sachin Kadyan, Lukas Jarosch, Dan Berenberg, Ian Fisk, Andrew Watkins, Stephen Ra, Richard Bonneau, and Mohammed AlQuraishi. Openproteinset: Training data for structural biology at scale. *Advances in Neural Information Processing Systems*, 36, 2024.
- [139] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [140] Shashikant Pujar, Nuala A O’Leary, Catherine M Farrell, Jane E Loveland, Jonathan M Mudge, Craig Wallin, Carlos G Girón, Mark Diekhans, If Barnes, Ruth Bennett, et al. Consensus coding sequence (ccds) database: a standardized set of human and mouse protein-coding regions supported by expert curation. *Nucleic Acids Research*, 46(D1):D221–D228, 2018.
- [141] Emmanuel Boutet, Damien Lieberherr, Michael Tognolli, Michel Schneider, and Amos Bairoch. Uniprotkb/swiss-prot: the manually annotated section of the uniprot knowledgebase. In *Plant bioinformatics: methods and protocols*, pages 89–112. Springer, 2007.
- [142] Baldwin Dumortier, Antoine Liutkus, Clément Carré, and Gabriel Krouk. Petribert: Augmenting bert with tridimensional encoding for inverse protein folding and design. *BioRxiv*, pages 2022–08, 2022.
- [143] Zuobai Zhang, Jiarui Lu, Vijil Chenthamarakshan, Aurélie Lozano, Payel Das, and Jian Tang. Structure-informed protein language model. *arXiv preprint arXiv:2402.05856*, 2024.
- [144] Sam Gelman, Bryce Johnson, Chase Freschlin, Sameer D’Costa, Anthony Gitter, and Philip A Romero. Biophysics-based protein language models for protein engineering. *bioRxiv*, pages 2024–03, 2024.
- [145] Rebecca F Alford, Andrew Leaver-Fay, Jeliazko R Jeliazkov, Matthew J O’Meara, Frank P DiMaio, Hahnbeom Park, Maxim V Shapovalov, P Douglas Renfrew, Vikram K Mulligan, Kalli Kappel, et al. The rosetta all-atom energy function for macromolecular modeling and design. *Journal of chemical theory and computation*, 13(6):3031–3048, 2017.
- [146] Roshan Rao, Joshua Meier, Tom Sercu, Sergey Ovchinnikov, and Alexander Rives. Transformer protein language models are unsupervised structure learners. *Biorxiv*, pages 2020–12, 2020.
- [147] Dexiong Chen, Philip Hartout, Paolo Pellizzoni, Carlos Oliver, and Karsten Borgwardt. Endowing protein language models with structural knowledge. *arXiv preprint arXiv:2401.14819*, 2024.
- [148] Zaixiang Zheng, Yifan Deng, Dongyu Xue, Yi Zhou, Fei Ye, and Quanquan Gu. Structure-informed language models are protein designers. In *International Conference on Machine Learning*, pages 42317–42338. PMLR, 2023.
- [149] Jeffrey A Ruffolo, Aadyot Bhatnagar, Joel Beazer, Stephen Nayfach, Jordan Russ, Emily Hill, Riffat Hussain, Joseph Gallagher, and Ali Madani. Adapting protein language models for structure-conditioned design. *bioRxiv*, pages 2024–08, 2024.
- [150] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- [151] Michel Van Kempen, Stephanie S Kim, Charlotte Tumescheit, Milot Mirdita, Jeongjae Lee, Cameron LM Gilchrist, Johannes Söding, and Martin Steinegger. Fast and accurate protein structure search with foldseek. *Nature Biotechnology*, 42(2):243–246, 2024.
- [152] Xiaohan Lin, Zhenyu Chen, Yanheng Li, Xingyu Lu, Chuanliu Fan, Ziqiang Cao, Shihao Feng, Yi Qin Gao, and Jun Zhang. Protokens: A machine-learned language for compact and informative encoding of protein 3d structures. *bioRxiv*, pages 2023–11, 2023.
- [153] Zhangyang Gao, Cheng Tan, and Stan Z Li. Foldtoken4: Consistent & hierarchical fold language. *bioRxiv*, pages 2024–08, 2024.
- [154] Barthelemy Meynard-Piganeau, Jiayou Zhang, James Gong, Xingyi Cheng, Yingtao Luo, Hugo Ly, Le Song, and Eric P Xing. Balancing locality and reconstruction in protein structure tokenizer. *bioRxiv*, pages 2024–12, 2024.
- [155] Michael Heinzinger, Konstantin Weissenow, Joaquin Gomez Sanchez, Adrian Henkel, Martin Steinegger, and Burkhard Rost. Prostr5: Bilingual language model for protein sequence and structure. *bioRxiv*, pages 2023–07, 2023.
- [156] Jin Su, Zhikai Li, Chenchen Han, Yuyang Zhou, Yan He, Junjie Shan, Xibin Zhou, Xing Chang, Dacheng Ma, OPMC, et al. Saprothub: Making protein modeling accessible to all biologists. *bioRxiv*, pages 2024–05, 2024.
- [157] Mingchen Li, Yang Tan, Bozitao Zhong, Ziyi Zhou, Huiqun Yu, Xinzhu Ma, Wanli Ouyang, Liang Hong, Bingxin Zhou, and Pan Tan. Deprot: A protein language model with quantized structure and disentangled attention. *bioRxiv*, pages 2024–04, 2024.

- [158] Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter Abbeel, and Yun Song. Evaluating protein transfer learning with tape. *Advances in neural information processing systems*, 32, 2019.
- [159] Minghao Xu, Zuobai Zhang, Jiarui Lu, Zhaocheng Zhu, Yangtian Zhang, Ma Chang, Runcheng Liu, and Jian Tang. Peer: a comprehensive and multi-task benchmark for protein sequence understanding. *Advances in Neural Information Processing Systems*, 35:35156–35173, 2022.
- [160] Pascal Notin, Aaron Kollasch, Daniel Ritter, Lood Van Niekerk, Steffanie Paul, Han Spinner, Nathan Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, et al. Proteingym: Large-scale benchmarks for protein fitness prediction and design. *Advances in Neural Information Processing Systems*, 36, 2024.
- [161] Xinyou Wang, Zaixiang Zheng, Fei Ye, Dongyu Xue, Shujian Huang, and Quanquan Gu. Dplm-2: A multimodal diffusion protein language model. *arXiv preprint arXiv:2410.13782*, 2024.
- [162] Nadav Brandes, Dan Ofer, Yam Peleg, Nadav Rappoport, and Michal Linial. Proteinbert: a universal deep-learning model of protein sequence and function. *Bioinformatics*, 38(8):2102–2110, 2022.
- [163] Zeyuan Wang, Qiang Zhang, HU Shuang-Wei, Haoran Yu, Xurui Jin, Zhichen Gong, and Huajun Chen. Multi-level protein structure pre-training via prompt learning. In *The Eleventh International Conference on Learning Representations*, 2022.
- [164] Geraldene Munsamy, Sebastian Lindner, Philipp Lorenz, and Noelia Ferruz. Zymctrl: a conditional language model for the controllable generation of artificial enzymes. In *NeurIPS Machine Learning in Structural Biology Workshop*, 2022.
- [165] Jannis Born and Matteo Manica. Regression transformer enables concurrent sequence regression and generation for molecular language modelling. *Nature Machine Intelligence*, 5(4):432–444, 2023.
- [166] Shengchao Liu, Yutao Zhu, Jiarui Lu, Zhao Xu, Weili Nie, Anthony Gitter, Chaowei Xiao, Jian Tang, Hongyu Guo, and Anima Anandkumar. A text-guided protein design framework. *arXiv preprint arXiv:2302.04611*, 2023.
- [167] Minghao Xu, Xinyu Yuan, Santiago Miret, and Jian Tang. Protst: Multi-modality learning of protein sequences and biomedical texts. In *International Conference on Machine Learning*, pages 38749–38767. PMLR, 2023.
- [168] Mingze Yin, Hanjing Zhou, Yiheng Zhu, Miao Lin, Yixuan Wu, Jialu Wu, Hongxia Xu, Chang-Yu Hsieh, Tingjun Hou, Jintai Chen, et al. Multi-modal clip-informed protein editing. *bioRxiv*, pages 2024-07, 2024.
- [169] Kevin E Wu, Howard Chang, and James Zou. Proteinclip: enhancing protein language models with natural language. *bioRxiv*, pages 2024-05, 2024.
- [170] Jin Su, Xibin Zhou, Xuting Zhang, and Fajie Yuan. Protrek: Navigating the protein universe through tri-modal contrastive learning. *bioRxiv*, pages 2024-05, 2024.
- [171] Ningyu Zhang, Zhen Bi, Xiaozhuan Liang, Siyuan Cheng, Haosen Hong, Shumin Deng, Jiazhang Lian, Qiang Zhang, and Huajun Chen. Ontoprotein: Protein pretraining with gene ontology embedding. *arXiv preprint arXiv:2201.11147*, 2022.
- [172] Hong-Yu Zhou, Yunxiang Fu, Zhicheng Zhang, Bian Cheng, and Yizhou Yu. Protein representation learning via knowledge enhanced primary structure reasoning. In *The Eleventh International Conference on Learning Representations*, 2022.
- [173] Young Su Ko, Jonathan Parkinson, and Wei Wang. Benchmarking text-integrated protein language model embeddings and embedding fusion on diverse downstream tasks. *bioRxiv*, pages 2024-08, 2024.
- [174] Lasse M Blaabjerg, Nicolas Jonsson, Wouter Boomsma, Amelie Stein, and Kresten Lindorff-Larsen. Ssemb: A joint embedding of protein sequence and structure enables robust variant effect predictions. *Nature Communications*, 15(1):9646, 2024.
- [175] Youhan Lee, Hasun Yu, Jaemyung Lee, and Jaehoon Kim. Pre-training sequence, structure, and surface features for comprehensive protein representation learning. In *The Twelfth International Conference on Learning Representations*, 2023.
- [176] Louis Robinson, Timothy Atkinson, Liviu Copoiu, Patrick Bordes, Thomas Pierrot, and Thomas D Barrett. Contrasting sequence with structure: Pre-training graph representations with plms. *bioRxiv*, pages 2023-12, 2023.
- [177] Fan Hu, Yishen Hu, Weihong Zhang, Huazhen Huang, Yi Pan, and Peng Yin. A multimodal protein representation framework for quantifying transferability across biochemical downstream tasks. *Advanced Science*, 10(22):2301223, 2023.
- [178] Tymor Hamamsy, Meet Barot, James T Morton, Martin Steinegger, Richard Bonneau, and Kyunghyun Cho. Learning sequence, structure, and function representations of proteins with language models. *bioRxiv*, 2023.
- [179] Mihaly Varadi, Stephen Anyango, Mandar Deshpande, Sreenath Nair, Cindy Natassia, Galabina Yordanova, David Yuan, Oana Stroe, Gemma Wood, Agata Laydon, et al. Alphafold protein structure database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic acids research*, 50(D1):D439–D444, 2022.
- [180] Jie Hou, Badri Adhikari, and Jianlin Cheng. Deepsf: deep convolutional neural network for mapping protein sequences to folds. *Bioinformatics*, 34(8):1295–1303, 2018.
- [181] Christine A Orengo, Alex D Michie, Susan Jones, David T Jones, Mark B Swindells, and Janet M Thornton. Cath—a hierarchic classification of protein domain structures. *Structure*, 5(8):1093–1109, 1997.
- [182] Typhaine Paysan-Lafosse, Matthias Blum, Sara Chuguransky, Tiago Grego, Beatriz Lázaro Pinto, Gustavo A Salazar, Maxwell L Bileschi, Peer Bork, Alan Bridge, Lucy Colwell, et al. Interpro in 2022. *Nucleic acids research*, 51(D1):D418–D427, 2023.
- [183] Philip Jones, David Binns, Hsin-Yu Chang, Matthew Fraser, Weizhong Li, Craig McAnulla, Hamish McWilliam, John Maslen, Alex Mitchell, Gift Nuka, et al. Interproscan 5: genome-scale protein function classification. *Bioinformatics*, 30(9):1236–1240, 2014.
- [184] Damian Szklarczyk, Annika L Gable, David Lyon, Alexander Junge, Stefan Wyder, Jaime Huerta-Cepas, Milan Simonovic, Nadezhda T Doncheva, John H Morris, Peer Bork, et al. String v11: protein–protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic acids research*, 47(D1):D607–D613, 2019.
- [185] Scott Federhen. The ncbi taxonomy database. *Nucleic acids research*, 40(D1):D136–D143, 2012.
- [186] Iz Beltagy, Kyle Lo, and Arman Cohan. Scibert: A pretrained language model for scientific text. *arXiv preprint arXiv:1903.10676*, 2019.
- [187] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23, 2021.
- [188] Zeyuan Wang, Qiang Zhang, Keyan Ding, Ming Qin, Xiang Zhuang, Xiaotong Li, and Huajun Chen. Instructprotein: Aligning human and protein language via knowledge instruction. *arXiv preprint arXiv:2310.03269*, 2023.
- [189] Le Zhuo, Zewen Chi, Minghao Xu, Heyan Huang, Heqi Zheng, Conghui He, Xian-Ling Mao, and Wentao Zhang. Protllm: An interleaved protein-language llm with protein-as-word pre-training. *arXiv preprint arXiv:2403.07920*, 2024.
- [190] Xibin Zhou, Chenchen Han, Yingqi Zhang, Jin Su, Kai Zhuang, Shiyu Jiang, Zichen Yuan, Wei Zheng, Fengyuan Dai, Yuyang Zhou, et al. Decoding the molecular language of proteins with evolla. *bioRxiv*, pages 2025-01, 2025.
- [191] Yuesen Li, Chengyi Gao, Xin Song, Xiangyu Wang, Yungang Xu, and Suxia Han. Druggpt: A gpt-based strategy for designing potential ligands targeting specific proteins. *bioRxiv*, pages 2023-06, 2023.
- [192] Kangjie Zheng, Siyu Long, Tianyu Lu, Junwei Yang, Xinyu Dai, Ming Zhang, Zaiqing Nie, Wei-Ying Ma, and Hao Zhou. Multi-scale protein language model for unified molecular modeling. *bioRxiv*, pages 2024-03, 2024.
- [193] Kathi Canese and Sarah Weis. Pubmed: the bibliographic database. *The NCBI handbook*, 2(1), 2013.
- [194] Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Huajun Chen. Mol-instructions: A large-scale biomolecular instruction dataset for large language models. *arXiv preprint arXiv:2306.08018*, 2023.
- [195] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [196] John J Irwin, Khanh G Tang, Jennifer Young, Chinzorig Dandarchuluun, Benjamin R Wong, Munkhzul Khurelbataar, Yurii S Moroz, John Mayfield, and Roger A Sayle. Zinc20—a free ultralarge-scale chemical database for ligand discovery. *Journal of chemical information and modeling*, 60(12):6065–6073, 2020.

- [197] Gengmo Zhou, Zhifeng Gao, Qiankun Ding, Hang Zheng, Hongteng Xu, Zhewei Wei, Linfeng Zhang, and Guolin Ke. Unimol: A universal 3d molecular representation learning framework. 2023.
- [198] Yizhen Luo, Jiahuan Zhang, Siqi Fan, Kai Yang, Yushuai Wu, Mu Qiao, and Zaiqing Nie. Biomedgpt: Open multimodal generative pre-trained transformer for biomedicine. *arXiv preprint arXiv:2308.09442*, 2023.
- [199] Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindulyte, Jia He, Siqian He, Qingliang Li, Benjamin A Shoemaker, Paul A Thiessen, Bo Yu, et al. Pubchem 2023 update. *Nucleic acids research*, 51(D1):D1373–D1380, 2023.
- [200] Shengchao Liu, Hanchen Wang, Weiyang Liu, Joan Lasenby, Hongyu Guo, and Jian Tang. Pre-training molecular graph representation with 3d geometry. *arXiv preprint arXiv:2110.07728*, 2021.
- [201] Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. *arXiv preprint arXiv:2402.17810*, 2024.
- [202] Xiang Zhuang, Keyan Ding, Tianwen Lyu, Yinuo Jiang, Xiaotong Li, Zhuoyi Xiang, Zeyuan Wang, Ming Qin, Kehua Feng, Jike Wang, et al. Instructbiomol: Advancing biomolecule understanding and design following human instructions. *arXiv preprint arXiv:2410.07919*, 2024.
- [203] Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke, Kyunghyun Cho, and Heng Ji. Translation between molecules and natural language. *arXiv preprint arXiv:2204.11817*, 2022.
- [204] Michael K Gilson, Tiqing Liu, Michael Baitaluk, George Nicola, Linda Hwang, and Jenny Chong. Bindingdb in 2015: a public database for medicinal chemistry, computational chemistry and systems pharmacology. *Nucleic acids research*, 44(D1):D1045–D1053, 2016.
- [205] Parit Bansal, Anne Morgat, Kristian B Axelsen, Venkatesh Muthukrishnan, Elisabeth Coudert, Lucila Aimò, Nevila Hyka-Nouspikel, Elisabeth Gasteiger, Arnaud Kerhornou, Teresa Batista Neto, et al. Rhea, the reaction knowledgebase in 2022. *Nucleic acids research*, 50(D1):D693–D700, 2022.
- [206] Hanwen Xu, Addie Woicik, Hoifung Poon, Russ B Altman, and Sheng Wang. Multilingual translation for zero-shot biomedical classification using biotranslator. *Nature Communications*, 14(1):738, 2023.
- [207] Sebastian Köhler, Michael Gargano, Nicolas Matentzoglou, Leigh C Carmody, David Lewis-Smith, Nicole A Vasilevsky, Daniel Danis, Ganna Balagura, Gareth Baynam, Amy M Brower, et al. The human phenotype ontology in 2021. *Nucleic acids research*, 49(D1):D1207–D1217, 2021.
- [208] Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*, 2022.
- [209] Zifeng Wang, Zichen Wang, Balasubramaniam Srinivasan, Vasilis N Ioannidis, Huzefa Rangwala, and Rishita Anubhai. Bio-bridge: Bridging biomedical foundation models via knowledge graph. *arXiv preprint arXiv:2310.03320*, 2023.
- [210] Gaetano T Montelione, Michael Nilges, Ad Bax, Peter Güntert, Torsten Herrmann, Jane S Richardson, Charles D Schwieters, Wim F Vranken, Geerten W Vuister, David S Wishart, et al. Recommendations of the wwPDB NMR validation task force. *Structure*, 21(9):1563–1570, 2013.
- [211] Guo-Wei Wei. Protein structure prediction beyond alphafold. *Nature Machine Intelligence*, 1(8):336–337, 2019.
- [212] Ruidong Wu, Fan Ding, Rui Wang, Rui Shen, Xiwen Zhang, Shitong Luo, Chenpeng Su, Zuofan Wu, Qi Xie, Bonnie Berger, et al. High-resolution de novo structure prediction from primary sequence. *BioRxiv*, pages 2022–07, 2022.
- [213] Ratul Chowdhury, Nazim Bouatta, Surojit Biswas, Christina Floristean, Anant Kharkar, Koushik Roy, Charlotte Rochereau, Gustaf Ahndritz, Joanna Zhang, George M Church, et al. Single-sequence protein structure prediction using a language model and deep learning. *Nature Biotechnology*, 40(11):1617–1623, 2022.
- [214] Jeffrey A Ruffolo, Lee-Shin Chu, Sai Pooja Mahajan, and Jeffrey J Gray. Fast, accurate antibody structure prediction from deep learning on massive set of natural antibodies. *Nature communications*, 14(1):2389, 2023.
- [215] Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, pages 1–3, 2024.
- [216] Uniprot: the universal protein knowledgebase in 2023. *Nucleic acids research*, 51(D1):D523–D531, 2023.
- [217] Yidong Song, Qianmu Yuan, Sheng Chen, Ken Chen, Yaoqi Zhou, and Yuedong Yang. Fast and accurate protein intrinsic disorder prediction by using a pretrained language model. *Briefings in bioinformatics*, 24(4):bbad173, 2023.
- [218] Kit Sang Chu and Justin B Siegel. Protein stability prediction by fine-tuning a protein language model on a mega-scale dataset. *bioRxiv*, pages 2023–11, 2023.
- [219] Raghvendra Mall, Ankita Singh, Chirag N Patel, Gregory Guirimand, and Filippo Castiglione. Vish-pred: an ensemble of fine-tuned esm models for protein toxicity prediction. *Briefings in Bioinformatics*, 25(4), 2024.
- [220] Boming Kang, Rui Fan, Chunmei Cui, and Qinghua Cui. Comprehensive prediction and analysis of human protein essentiality based on a pretrained large language model. *Nature Computational Science*, pages 1–11, 2024.
- [221] Chakradhar Guntuboina, Adrita Das, Parisa Mollaei, Seongwon Kim, and Amir Barati Farimani. Peptidebert: A language model based on transformers for peptide property prediction. *The Journal of Physical Chemistry Letters*, 14(46):10427–10434, 2023.
- [222] Xuenan Mi, Susanna E Barrett, Douglas A Mitchell, and Diwakar Shukla. Lassoesm: A tailored language model for enhanced lasso peptide property prediction. *bioRxiv*, pages 2024–10, 2024.
- [223] Mark Zaretskii, Pavel Buslaev, Igor Kozlovskii, Alexander Morozov, and Petr Popov. Approaching optimal pH enzyme prediction with large language models. *ACS Synthetic Biology*, 13(9):3013–3021, 2024.
- [224] Japheth E Gado, Matthew Knotts, Ada Y Shaw, Debora Marks, Nicholas P Gauthier, Chris Sander, and Gregg T Beckham. Deep learning prediction of enzyme optimum pH. *bioRxiv*, pages 2023–06, 2023.
- [225] Nadav Brandes, Grant Goldman, Charlotte H Wang, Chun Jimmie Ye, and Vasilis Ntranos. Genome-wide prediction of disease variant effects with a deep protein language model. *Nature Genetics*, 55(9):1512–1522, 2023.
- [226] Ziyi Zhou, Liang Zhang, Yuanxi Yu, Banghao Wu, Mingchen Li, Liang Hong, and Pan Tan. Enhancing efficiency of protein language models with minimal wet-lab data through few-shot learning. *Nature Communications*, 15(1):5566, 2024.
- [227] Zeyuan Wang, Keyan Ding, Ming Qin, Xiaotong Li, Xiang Zhuang, Yu Zhao, Jianhua Yao, Qiang Zhang, and Huajun Chen. Deplm: Denoising protein language models for property optimization. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- [228] Zachary N Flamholz, Steven J Biller, and Libusha Kelly. Large language models improve annotation of prokaryotic viral proteins. *Nature Microbiology*, pages 1–13, 2024.
- [229] Yidong Song, Qianmu Yuan, Sheng Chen, Yuansong Zeng, Huiying Zhao, and Yuedong Yang. Accurately predicting enzyme functions through geometric graph learning on esmfold-predicted structures. *Nature Communications*, 15(1):8180, 2024.
- [230] Shaojun Wang, Ronghui You, Yunjia Liu, Yi Xiong, and Shanfeng Zhu. Netgo 3.0: protein language model improves large-scale functional annotations. *Genomics, Proteomics & Bioinformatics*, 21(2):349–358, 2023.
- [231] Maxat Kulmanov, Francisco J Guzmán-Vega, Paula Duek Roggli, Lydie Lane, Stefan T Arold, and Robert Hoehndorf. Protein function prediction as approximate semantic entailment. *Nature Machine Intelligence*, pages 1–9, 2024.
- [232] Qianmu Yuan, Junjie Xie, Jiancong Xie, Huiying Zhao, and Yuedong Yang. Fast and accurate protein function prediction from sequence through pretrained language model and homology-based label diffusion. *Briefings in bioinformatics*, 24(3):bbad117, 2023.
- [233] Yuansong Zeng, Zhuoyi Wei, Qianmu Yuan, Sheng Chen, Weijiang Yu, Yutong Lu, Jianzhao Gao, and Yuedong Yang. Identifying b-cell epitopes using alphafold2 predicted structures and pretrained language model. *Bioinformatics*, 39(4):btad187, 2023.
- [234] Felix Teufel, José Juan Almagro Armenteros, Alexander Rosenberg Johansen, Magnús Halldór Gíslason, Silas Irby Pihl, Konstantinos D Tsigiris, Ole Winther, Søren Brunak, Gunnar von Heijne, and Henrik Nielsen. Signalp 6.0 predicts all five types of signal peptides using protein language models. *Nature biotechnology*, 40(7):1023–1025, 2022.

- [235] Junbo Shen, Qinqin Yu, Shenyang Chen, Qingxiong Tan, Jingchen Li, and Yu Li. Unbiased organism-agnostic and highly sensitive signal peptide predictor with deep protein language model. *Nature Computational Science*, 4(1):29–42, 2024.
- [236] Hilbert Yuen In Lam, Jia Sheng Guan, Xing Er Ong, Robbe Pincket, and Yuguang Mu. Protein language models are performant in structure-free virtual screening. *Briefings in Bioinformatics*, 25(6):bbae480, 2024.
- [237] Shion Honda, Shoi Shi, and Hiroki R Ueda. Smiles transformer: Pre-trained molecular fingerprint for low data drug discovery. *arXiv preprint arXiv:1911.04738*, 2019.
- [238] Rohan Gorantla, Aryo Pradipta Gema, Ian Xi Yang, Álvaro Serrano-Morrás, Benjamin Suutari, Jordi Juárez Jiménez, and Antonia SJS Mey. Learning binding affinities via fine-tuning of protein and ligand language models. *bioRxiv*, pages 2024–11, 2024.
- [239] Walid Ahmad, Elana Simon, Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. Chemberta-2: Towards chemical foundation models. *arXiv preprint arXiv:2209.01712*, 2022.
- [240] Chenqing Hua, Bozita Zhong, Sitao Luan, Liang Hong, Guy Wolf, Doina Precup, and Shuangjia Zheng. Reactzyme: A benchmark for enzyme-reaction prediction. *arXiv preprint arXiv:2408.13659*, 2024.
- [241] Xiaorui Wang, Xiaodan Yin, Dejun Jiang, Huifeng Zhao, Zhenxing Wu, Odin Zhang, Jike Wang, Yuquan Li, Yafeng Deng, Huanxiang Liu, et al. Multi-modal deep learning enables efficient and accurate annotation of enzymatic active sites. *Nature Communications*, 15(1):7348, 2024.
- [242] Zuobai Zhang, Minghao Xu, Aurelie Lozano, Vijil Chenthamarakshan, Payel Das, and Jian Tang. Enhancing protein language model with structure-based encoder and pre-training. In *ICLR 2023-Machine Learning for Drug Discovery workshop*, 2023.
- [243] Garyk Brixi, Tianzheng Ye, Lauren Hong, Tian Wang, Connor Monticello, Natalia Lopez-Barbosa, Sophia Vincoff, Vivian Yuditstra, Lin Zhao, Elena Haarer, et al. Salt&peppr is an interface-predicting language model for designing peptide-guided protein degraders. *Communications Biology*, 6(1):1081, 2023.
- [244] Samuel Sledzieski, Meghana Kshirsagar, Minkyung Baek, Rahul Dodhia, Juan Lavista Ferras, and Bonnie Berger. Democratizing protein language models with parameter-efficient fine-tuning. *Proceedings of the National Academy of Sciences*, 121(26):e2405840121, 2024.
- [245] Mingyu Jin, Xue Haochen, Zhenting Wang, Boming Kang, Ruosong Ye, Kaixiong Zhou, Mengnan Du, and Yongfeng Zhang. Prollm: Protein chain-of-thoughts enhanced llm for protein-protein interaction prediction. *bioRxiv*, pages 2024–04, 2024.
- [246] Mahdi Pourmirzaei, Farzaneh Esmaili, Mohammadreza Pourmirzaei, Duolin Wang, and Dong Xu. Prot2token: A multi-task framework for protein language processing using autoregressive language modeling. *bioRxiv*, pages 2024–05, 2024.
- [247] Gayane Chilingaryan, Hovhannes Tamoyan, Ani Tevosyan, Nelly Babayan, Lusine Khondkaryan, Karen Hambardzumyan, Zaven Navoyan, Hrant Khachatryan, and Armen Aghajanyan. Bartsmls: Generative masked language models for molecular representations. *arXiv preprint arXiv:2211.16349*, 2022.
- [248] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2(3):6, 2023.
- [249] Yijia Xiao, Edward Sun, Yiqiao Jin, Qifan Wang, and Wei Wang. Proteingpt: Multimodal llm for protein property prediction and structure understanding. *arXiv preprint arXiv:2408.11363*, 2024.
- [250] Wenkai Xiang, Zhaoping Xiong, Huan Chen, Jiacheng Xiong, Wei Zhang, Zunyun Fu, Mingyue Zheng, Bing Liu, and Qian Shi. Fapm: functional annotation of proteins using multimodal models beyond structural modeling. *Bioinformatics*, 40(12):btae680, 2024.
- [251] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [252] Zhiyuan Liu, An Zhang, Hao Fei, Enzhi Zhang, Xiang Wang, Kenji Kawaguchi, and Tat-Seng Chua. Prott3: Protein-to-text generation for text-based protein understanding. *arXiv preprint arXiv:2405.12564*, 2024.
- [253] Frances Ding and Jacob Noah Steinhardt. Protein language models are biased by unequal sequence sampling across the tree of life. *bioRxiv*, pages 2024–03, 2024.
- [254] Cade Gordon, Amy X Lu, and Pieter Abbeel. Protein language model fitness is a matter of preference. *bioRxiv*, pages 2024–10, 2024.
- [255] Ronghui You, Shuwei Yao, Yi Xiong, Xiaodi Huang, Fengzhu Sun, Hiroshi Mamitsuka, and Shanfeng Zhu. Netgo: improving large-scale protein function prediction with massive network information. *Nucleic acids research*, 47(W1):W379–W387, 2019.
- [256] Shuwei Yao, Ronghui You, Shaojun Wang, Yi Xiong, Xiaodi Huang, and Shanfeng Zhu. Netgo 2.0: improving large-scale protein function prediction with massive sequence, text, domain, family and network information. *Nucleic acids research*, 49(W1):W469–W475, 2021.
- [257] Frances H Arnold. Directed evolution: bringing new chemistry to life. *Angewandte Chemie (International Ed. in English)*, 57(16):4143, 2018.
- [258] Michael S Packer and David R Liu. Methods for the directed evolution of proteins. *Nature Reviews Genetics*, 16(7):379–394, 2015.
- [259] Brian L Hie, Varun R Shanker, Duo Xu, Theodora UJ Bruun, Payton A Weidenbacher, Shaogeng Tang, Wesley Wu, John E Pak, and Peter S Kim. Efficient evolution of human antibodies from general protein language models. *Nature Biotechnology*, 42(2):275–283, 2024.
- [260] Sean R Johnson, Sarah Monaco, Kenneth Massie, and Zaid Syed. Generating novel protein sequences using gibbs sampling of masked language models. *bioRxiv*, pages 2021–01, 2021.
- [261] Damiano Sgarbosa, Umberto Lupo, and Anne-Florence Bitbol. Generative power of a protein language model trained on multiple sequence alignments. *Elife*, 12:e79854, 2023.
- [262] Thanh VT Tran and Truong Son Hy. Protein design by directed evolution guided by large language models. *bioRxiv*, pages 2023–11, 2023.
- [263] Yves Gaetan Nana Teukam, Federico Zipoli, Teodoro Laino, Emanuele Criscuolo, Francesca Grisoni, and Matteo Manica. Integrating genetic algorithms and language models for enhanced enzyme design. 2024.
- [264] Hideki Yamaguchi and Yutaka Saito. Evoopt: an msa-guided, fully unsupervised sequence optimization pipeline for protein design. In *Machine Learning for Structural Biology Workshop, NeurIPS*, 2022.
- [265] Fan Jiang, Mingchen Li, Jiajun Dong, Yuanxi Yu, Xinyu Sun, Banghao Wu, Jin Huang, Liqi Kang, Yufeng Pei, Liang Zhang, et al. A general temperature-guided language model to design proteins of enhanced stability and activity. *Science Advances*, 10(48):eadr2641, 2024.
- [266] Jeremie Theddy Darmawan, Yarin Gal, and Pascal Notin. Sampling protein language models for functional protein design. In *NeurIPS 2023 Generative AI and Biology (GenBio) Workshop*, 2023.
- [267] Kaiyi Jiang, Zhaoqing Yan, Matteo Di Bernardo, Samantha R Sgrizzi, Lukas Villiger, Alisan Kayabolen, BJ Kim, Josephine K Carscadden, Masahiro Hiraizumi, Hiroshi Nishimasu, et al. Rapid in silico directed evolution by a protein language model with evolvepro. *Science*, page eadr6006, 2024.
- [268] Shengchao Liu, Jiong Xiao Wang, Yijin Yang, Chengpeng Wang, Ling Liu, Hongyu Guo, and Chaowei Xiao. Chatgpt-powered conversational drug editing using retrieval and domain feedback. *arXiv preprint arXiv:2305.18090*, 2023.
- [269] Alexander E Chu, Tianyu Lu, and Po-Ssu Huang. Sparks of function by de novo protein design. *Nature biotechnology*, 42(2):203–215, 2024.
- [270] Jiezhong Qiu, Junde Xu, Jie Hu, Hanqun Cao, Liya Hou, Zijun Gao, Xinyi Zhou, Anni Li, Xiujuan Li, Bin Cui, et al. Instructplm: Aligning protein language models to follow protein structure instructions. *bioRxiv*, pages 2024–04, 2024.
- [271] Fengyuan Dai, Yuliang Fan, Jin Su, Chentong Wang, Chenchen Han, Xibin Zhou, Jianming Liu, Hui Qian, Shunzhi Wang, Anping Zeng, et al. Toward de novo protein design from natural language. *bioRxiv*, pages 2024–08, 2024.
- [272] Geraldene Munsamy, Ramiro Illanes-Vicioso, Silvia Funcillo, Ioanna T Nakou, Sebastian Lindner, Gavin Ayres, Lesley S Sheehan, Steven Moss, Ulrich Eckhard, Philipp Lorenz, et al. Conditional language models enable the efficient design of proficient enzymes. *bioRxiv*, pages 2024–05, 2024.
- [273] Alireza Ghafarollahi and Markus J Buehler. Protagents: Protein discovery via large language model multi-agent collaborations combining physics and machine learning. *arXiv preprint arXiv:2402.04268*, 2024.
- [274] Klaus Rajewsky. Clonal selection and learning in the antibody system. *Nature*, 381(6585):751–758, 1996.



- [275] Andreas H Laustsen, Victor Greiff, Aneesh Karatt-Vellatt, Serge Muyldermans, and Timothy P Jenkins. Animal immunization, in vitro display technologies, and machine learning for antibody discovery. *Trends in Biotechnology*, 39(12):1263–1273, 2021.
- [276] Haohuai He, Bing He, Lei Guan, Yu Zhao, Feng Jiang, Guanxing Chen, Qingge Zhu, Calvin Yu-Chian Chen, Ting Li, and Jianhua Yao. De novo generation of sars-cov-2 antibody cdrh3 with a pre-trained generative large language model. *Nature Communications*, 15(1):6867, 2024.
- [277] Varun R Shanker, Theodora UJ Bruun, Brian L Hie, and Peter S Kim. Unsupervised evolution of protein and antibody complexes with a structure-informed language model. *Science*, 385(6704):46–53, 2024.
- [278] Peter K Robinson. Enzymes: principles and biotechnological applications. *Essays in biochemistry*, 59:1, 2015.
- [279] Philip A Romero, Tuan M Tran, and Adam R Abate. Dissecting enzyme function with microfluidic-based deep mutational scanning. *Proceedings of the National Academy of Sciences*, 112(23):7159–7164, 2015.
- [280] Roy Saffhill and JJ Weiss. Effects of ionizing radiation on in vitro replication of dna by dna polymerase i. *Nature New Biology*, 241(107):69–71, 1973.
- [281] Charles S Craik, Michael J Page, and Edwin L Madison. Proteases as therapeutics. *Biochemical Journal*, 435(1):1–16, 2011.
- [282] Uwe T Bornscheuer, GW Huisman, RJ Kazlauskas, S Lutz, JC Moore, and K Robins. Engineering the third wave of biocatalysis. *Nature*, 485(7397):185–194, 2012.
- [283] Sean R Johnson, Xiaozhi Fu, Sandra Viknander, Clara Goldin, Sarah Monaco, Aleksej Zelezniak, and Kevin K Yang. Computational scoring and experimental evaluation of enzymes generated by neural networks. *Nature biotechnology*, pages 1–10, 2024.
- [284] Fátima Gebauer and Matthias W Hentze. Molecular mechanisms of translational control. *Nature reviews Molecular cell biology*, 5(10):827–835, 2004.
- [285] Flora Cozzolino, Ilaria Iacobucci, Vittoria Monaco, and Maria Monti. Protein–dna/rna interactions: an overview of investigation methods in the-omics era. *Journal of Proteome Research*, 20(6):3018–3030, 2021.
- [286] Yutai Hou, Yingce Xia, Lijun Wu, Shufang Xie, Yang Fan, Jinhua Zhu, Tao Qin, and Tie-Yan Liu. Discovering drug–target interaction knowledge from biomedical literature. *Bioinformatics*, 38(22):5100–5107, 2022.
- [287] James P Hughes, Stephen Rees, S Barrett Kalindjian, and Karen L Philpott. Principles of early drug discovery. *British journal of pharmacology*, 162(6):1239–1249, 2011.
- [288] Yogesh Kalakoti, Shashank Yadav, and Durai Sundar. Transdti: transformer-based language models for estimating dtis and building a drug recommendation workflow. *ACS omega*, 7(3):2706–2717, 2022.
- [289] Henry Kenlay, Frédéric A Dreyer, Aleksandr Kovaltsuk, Dom Miketa, Douglas Pires, and Charlotte M Deane. Large scale paired antibody language models. *PLOS Computational Biology*, 20(12):e1012646, 2024.
- [290] Shuxian Zou, Hui Li, Shentong Mo, Xingyi Cheng, Eric Xing, and Le Song. Linker-tuning: Optimizing continuous prompts for heterodimeric protein prediction. *arXiv preprint arXiv:2312.01186*, 2023.
- [291] Shitong Luo, Yufeng Su, Zuofan Wu, Chenpeng Su, Jian Peng, and Jianzhu Ma. Rotamer density estimator is an unsupervised learner of the effect of mutations on protein-protein interaction. *bioRxiv*, pages 2023–02, 2023.
- [292] Anne Doerr. Protein design: the experts speak. *Nature Biotechnology*, 42(2):175–178, 2024.
- [293] Zhidian Zhang, Hannah K Wayment-Steele, Garyk Bixi, Haobo Wang, Dorothee Kern, and Sergey Ovchinnikov. Protein language models learn evolutionary statistics of interacting sequence motifs. *Proceedings of the National Academy of Sciences*, 121(45):e2406285121, 2024.
- [294] Elana Simon and James Zou. Interplm: Discovering interpretable features in protein language models via sparse autoencoders. *bioRxiv*, pages 2024–11, 2024.
- [295] Haoyi Xiong, Xiaofei Zhang, Jiamin Chen, Xinhao Sun, Yuchen Li, Zeyi Sun, Mengnan Du, et al. Towards explainable artificial intelligence (xai): A data mining perspective. *arXiv preprint arXiv:2401.04374*, 2024.
- [296] Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak, Shengchao Liu, Peter Van Katwyk, Andreea Deac, et al. Scientific discovery in the age of artificial intelligence. *Nature*, 620(7972):47–60, 2023.
- [297] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [298] Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. Augmenting large language models with chemistry tools. *Nature Machine Intelligence*, pages 1–11, 2024.
- [299] Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research with large language models. *Nature*, 624(7992):570–578, 2023.
- [300] Yaiza Serrano, Álvaro Ciudad, and Alexis Molina. Are protein language models compute optimal? *arXiv preprint arXiv:2406.07249*, 2024.
- [301] Xingyi Cheng, Bo Chen, Pan Li, Jing Gong, Jie Tang, and Le Song. Training compute-optimal protein language models. *bioRxiv*, pages 2024–06, 2024.
- [302] Quentin Fournier, Robert M Vernon, Almer van der Sloot, Benjamin Schulz, Sarath Chandar, and Christopher James Langmead. Protein language models: Is scaling necessary? *bioRxiv*, pages 2024–09, 2024.
- [303] Luiz C Vieira, Morgan C Handojo, and Claus O Wilke. Scaling down for efficiency: Medium-sized transformer models for protein sequence transfer learning. *bioRxiv*, pages 2024–11, 2024.
- [304] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022.
- [305] Muhammed Hasan Celik and Xiaohui Xie. Efficient inference, training, and fine-tuning of protein language models. *bioRxiv*, pages 2024–10, 2024.
- [306] Fred Zhangzhi Peng, Pranam Chatterjee, and contributors. Faesm: An efficient pytorch implementation of evolutionary scale modeling (esm). <https://github.com/pengzhangzhi/faesm>, 2024. Efficient PyTorch implementation of ESM with FlashAttention and Scalar Dot-Product Attention (SDPA).