

Deep Learning Methods for Protein Function Prediction

Frimpong Boadu¹, Ahhyun Lee¹, and Jianlin Cheng¹

¹University of Missouri

March 12, 2024

Abstract

Predicting protein function from protein sequence, structure, interaction, and other relevant information is important for generating hypotheses for biological experiments and studying biological systems, and therefore has been a major challenge in protein bioinformatics. Numerous computational methods had been developed to advance protein function prediction gradually in the last two decades. Particularly, in the recent years, leveraging the revolutionary advances in artificial intelligence (AI), more and more deep learning methods have been developed to improve protein function prediction at a faster pace. Here, we provide an in-depth review of the recent developments of deep learning methods for protein function prediction. We summarize the significant advances in the field, identify several remaining major challenges to be tackled, and suggest some potential directions to explore. The data sources and evaluation metrics widely used in protein function prediction are also discussed to assist the machine learning, AI, and bioinformatics communities to develop more cutting-edge methods to advance protein function prediction.

PROTEIN FUNCTION PREDICTION

Deep Learning Methods for Protein Function Prediction

Frimpong Boadu^{1, †}, Ahhyun Lee^{1, †} and Jianlin Cheng^{1, *}¹Department of Electrical Engineering and Computer Science, University of Missouri, Columbia, 65211, MO, USA[†]These authors contributed equally. *Corresponding author: chengji@missouri.edu

FOR PUBLISHER ONLY Received on Date Month Year; revised on Date Month Year; accepted on Date Month Year

Abstract

Predicting protein function from protein sequence, structure, interaction, and other relevant information is important for generating hypotheses for biological experiments and studying biological systems, and therefore has been a major challenge in protein bioinformatics. Numerous computational methods had been developed to advance protein function prediction gradually in the last two decades. Particularly, in the recent years, leveraging the revolutionary advances in artificial intelligence (AI), more and more deep learning methods have been developed to improve protein function prediction at a faster pace. Here, we provide an in-depth review of the recent developments of deep learning methods for protein function prediction. We summarize the significant advances in the field, identify several remaining major challenges to be tackled, and suggest some potential directions to explore. The data sources and evaluation metrics widely used in protein function prediction are also discussed to assist the machine learning, AI, and bioinformatics communities to develop more cutting-edge methods to advance protein function prediction.

Key words: protein function prediction, deep learning, artificial intelligence, gene ontology

1. Introduction

Proteins are essential molecules in all living organisms. Their role encompasses structural support, biochemical catalysis, gene regulation, enzymatic activities, and signal transduction[1, 2]. Determining the functions of proteins is a key step to understand biological systems and modulate biological processes, which plays an important role in biomedical research and biotechnology development. Furthermore, proteins are common targets in drug discovery[3, 4, 5] because many proteins are implicated in diseases, and protein function information can facilitate the development of drugs targeting them. As the structure of protein can be determined by experimental techniques such as x-ray crystallography, the function of proteins can also be determined by experimental techniques such as biochemical assays and enzymatic analysis. However, the experimental techniques for protein function determination is expensive, time-consuming, and labor-intensive and can only be applied to a small number of proteins. Therefore, making precise protein function prediction computationally holds the key to address the need of function information for most proteins and has become a critical challenge in bioinformatics.

Currently, hundreds of millions of protein sequences have been generated through numerous genome and transcriptome sequencing projects. However, less than 1% of them have experimentally determined protein function information. This presents a huge gap between known protein sequences and

their functions. Therefore, it is critical to devise advanced computational methods to accurately predict protein function to fill the gap as the recent development of deep learning methods has done for protein structure prediction [6, 7, 8, 9].

A plethora of various computational methods have been developed to predict protein function, many of which had been reviewed and assessed previously[10, 11, 12]. Recently, as AI is transforming many scientific fields, cutting-edge prediction methods based on deep learning approaches have been thriving in the protein function prediction field, leading to a significant improvement of prediction accuracy over the previous generation of computational protein function prediction methods. Therefore, there is a need of reviewing these latest advances to facilitate the development of more deep learning methods to address the remaining challenges in the field.

Here, we present a comprehensive overview of recent deep learning methods developed to advance protein function prediction. **Fig. 1** illustrates a general workflow of deep learning-based prediction of protein function defined by the gene ontology (GO) terms[13]. We classify these methods roughly into four main categories based on the input information used by them: (1) sequence-based methods of using only protein sequence as input, (2) structure-based methods of using protein structure as input, (3) Interaction-based methods of using protein-protein interaction information as input, and (4) integrative methods that using multiple sources of information

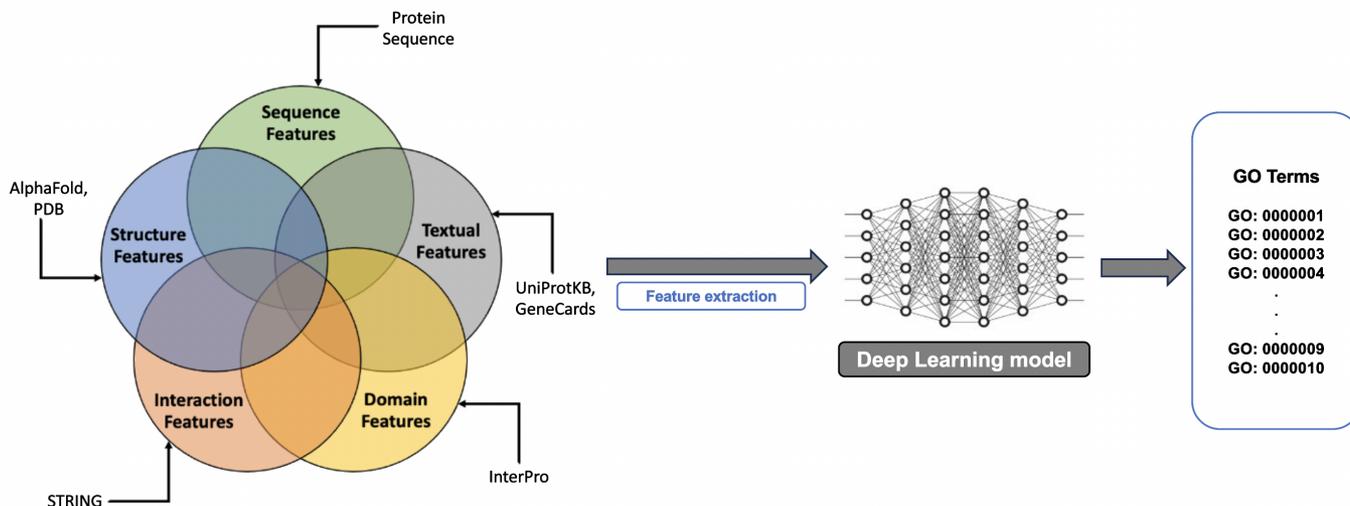


Fig. 1. The general workflow of deep learning-based protein function prediction. One or multiple sources of data such as protein sequences, protein structures (e.g., structures retrieved from the AlphaFoldDB[16] and the Protein Data Bank (PDB)[17]), protein-protein interaction from the STRING database[18], protein family and domain information from the Interpro database[19], and the textual description of proteins in the literature such as UniProt Knowledgebase(UniProtKB) [20] and GeneCards [21] are presented as input. The features are then extracted from the input data, which are fed into deep learning models to predict protein function as output. Protein function are usually described as gene ontology (GO)[13] function terms. Therefore, protein function prediction is essentially a classification problem. Because one protein may have multiple functions described by multiple GO terms, it is a multi-label classification problem.

as input. It is worth noting that structure-based or interaction-based methods often also use sequence information implicitly in addition to using structure or interaction information, but they are not classified as integrative methods. Moreover, we also discuss the latest few-shot learning [14, 15] paradigm that improves the prediction of rarely annotated protein function terms associated with few proteins. **Table 1** lists the types, input features, neural network architectures, and availability of 30 deep learning protein function prediction methods reviewed in this article. Furthermore, in addition to surveying the deep learning methods, we discuss the data sources, standard benchmarks (i.e., the Critical Assessment of Protein Function Annotation (CAFA)[10]), and evaluation metrics widely used for protein function prediction to assist the AI, machine learning, and bioinformatics communities to find necessary resources to develop more protein function prediction methods. Moreover, we identify several major remaining challenges in protein function prediction and envision that developing Large Language Models for Proteins (LLMP), akin to the Large Language Models (LLM) used in natural language processing (NLP), such as ChatGPT[14], can be a promising approach to addressing the challenges. These topics are discussed in detail in the sections below.

2. Sequence-based protein function prediction

Sequence-based prediction methods uses different kinds of deep learning architectures to take protein sequence information as input to predict protein function. Several deep learning models that have demonstrated effectiveness for dealing with sequential data are: 1) convolutional neural networks(CNNs) [22], 2) recurrent neural networks(RNNs) [23, 24], 3) deep neural networks(DNNs) [25, 26], and 4) attention-based transformers [2, 27]. CNNs are effective at identifying motifs (short conserved sequence patterns associated with distinct

protein functions), local patterns, and spatial relationships in the protein sequences. RNNs, particularly, Long Short-Term Memory networks(LSTMs) [28], can capture sequential dependence between amino acids in protein sequences. DNNs also hold significance in capturing the non-linear relationships between protein function and sequences through multiple neural network layers. Finally, the attention mechanism and transformer architecture are well-suited for sequence-based function prediction due to their ability to capture long-range dependencies between amino acids in protein sequences. Besides directly applying transformer based architectures to protein function prediction, several methods [29, 30, 31] leverage transformer-based pre-trained protein language models to extract representative embeddings from protein sequences for downstream protein function prediction tasks. In the subsequent sections below, we discuss the specific methods that harness these deep learning models to address the intricacies of predicting protein function from sequences.

2.1. RNN-based protein function prediction

ProLanGO[32] treats the protein function prediction problem as a language translation problem and applies a RNN-based Neural Machine Translation (NMT) model to tackle it. Protein sequences (input) and Gene Ontology terms (output) are regarded as two separate languages, ProLan and GoLan, respectively. Protein sequences are represented as a series of k-mers (i.e., a substring or word of k amino acids). Protein words are extracted based on the frequency of k-mers. GO function terms are generally represented as a directed acyclic tree structure based on their relationships, with each term uniquely identified by a seven-digit number. ProLanGo allows capturing the hierarchical relationship between GO terms and enables the sequence to function translation through the depth-first search(DFS). Each GO term is assigned to a 26-base Alphabet ID according to its order of being visited during

Table 1. The classification of 30 deep learning protein function prediction methods and their input features, architectures, and availability. Sequence, structure, interaction, and domain refers to four types of typical input features: sequence-based features, structure-based features, protein interaction-based features, and other features based on protein family and domain information. RNN stands for both standard recurrent neural networks and advanced ones like Gated Recurrent Unit (GRU) and Long Short-Term Memory (LSTM), CNN for convolutional neural networks, and GNN for graph neural networks. Attention denotes the methods utilizing self-attention mechanisms, transformers, and techniques extracting features from pre-trained attention- or transformer-based architectures. DNN refers to deep neural networks that use multilayer perceptrons (MLP) as a main part of the model architecture beyond using them in the final classification layer. Few-shot refers to methods specifically designed to utilize deep learning models for predicting GO terms with few annotations. We also include a link to the GitHub repository or webpage of the tool. For tools whose link we cannot find, we use NA.

Methods	Features					Deep Learning Architecture					Few-shot	URL
	Sequence	Structure	Interaction	Domain	Text	DNN	CNN	GNN	RNN	Attention		
ProLanGO	✓								✓			NA
FUTUSA	✓						✓					GitHub
DeepGOPlus	✓						✓					Web
PFmulDL	✓						✓		✓			GitHub
DEEPred	✓					✓						GitHub
TALE	✓									✓	✓	GitHub
TEMPROT	✓					✓				✓		GitHub
SPROF-GO	✓					✓				✓		GitHub Web
ATGO	✓					✓				✓		Web
PANDA2	✓							✓		✓		Web
DeepFRI	✓	✓						✓	✓			GitHub Web
GAT-GO	✓	✓					✓	✓		✓		GitHub
TransFun	✓	✓						✓		✓		GitHub
Struct2GO	✓	✓						✓		✓		GitHub
Mashup			✓									Web
deepNF			✓			✓						GitHub
MELISSA			✓									GitHub
NetQuilt	✓		✓			✓						GitHub
DeepGO	✓		✓			✓	✓					GitHub Web
STRING2GO			✓			✓						GitHub
DeepGraphGO			✓	✓				✓				GitHub
GRAPH2GO	✓		✓	✓				✓				GitHub
NetGO2	✓			✓	✓				✓			Web
NetGO3	✓			✓	✓					✓		Web
SDN2GO	✓		✓	✓		✓	✓					GitHub
PFP-GO	✓		✓	✓								Web
MultiPredGO	✓	✓	✓			✓	✓					GitHub
DeepGATGO	✓				✓	✓		✓		✓		NA
ProTranslator	✓		✓		✓		✓				✓	GitHub
DeepGOZero				✓		✓					✓	GitHub

Sequence-based

Structure-based

Interaction-based

Integrative

Few-Shot

the DFS traversal. Given the ProLan and GOLan languages, an encoder-decoder based on RNNs is trained to predict GOLan from ProLan. The encoder is used to encode a ProLan sentence into fixed-length vectors, and the decoder decodes the representation into a GOLan sentence. The network is

trained by maximizing the conditional probability of predicting a GOLan sentence given a ProLan sentence.

2.2. CNN-based protein function prediction

FUTUSA[33] has four components: CNN-based embedding layers, CNN-based feature extraction, dense layers, and a classification layer. The embedding layers are used to convert protein sequences to numerical vectors. To alleviate the limitations of one-hot encoding such as the inability to capture physiochemical properties of amino acids, a one-dimensional CNN is employed to generate the amino acid embedding vector, followed by another CNN to extract spatial features, whose output is fed into dense layers to generate hidden features. The hidden features are used by the final classification layer to predict GO terms.

DeepGOPlus[34] combines the function prediction from a CNN network and the sequence similarity to improve prediction accuracy. It uses one-dimensional CNN filters to learn similar patterns (motifs) in sequences. An input sequence is transformed into a matrix representation of dimension 21×2000 using a one-hot encoding strategy, where a one-hot vector of 21 binary numbers is used to represent an amino acid and the maximum number of amino acids to be represented is 2000. The input is fed into a set of CNN layers with varying filter sizes to generate features capturing sequence motifs of different size. The features are pooled together and selected by a MaxPooling layer. The output of the MaxPooling layer is forwarded to a fully connected classification layer to predict GO terms. DeepGOPlus is a general sequence-based protein function prediction that can be applied to proteins in any taxa or kingdom of species.

PFmulDL[35] integrates both a multi-kernel convolutional neural network and a gated recurrent unit (GRU) to predict protein function. Like DeepGoPlus, it employs a one-hot strategy to encode an input protein sequence. The encoding serves as input for a multi-kernel CNN model, which is fine-tuned by a pre-training process. The output layer of the CNN is used as input for the GRU to generate features, which are used as input for a fully connected layer to predict GO terms. In order to prevent issues such as gradient vanishing and overfitting, it uses transfer learning (TL) to improve training, leading to the improved performance of protein function prediction. Particularly, it enhances the prediction accuracy for 'rare GO terms (minority class)' without compromising the performance for the 'common GO terms (major classes)'.

2.3. DNN-based protein function prediction

DEEPred[36] employs a deep learning model organized as a stack of multi-task feed-forward deep neural networks (DNNs). Each DNN is independently designed to predict groups of 4 or 5 Gene Ontology (GO) terms. The grouping is based on the levels of GO terms in the GO graph, determined through the topological sorting. Groups are carefully created to ensure that GO terms within the same group have similar numbers of annotations, addressing the variability in protein associations. This approach aims to enhance the model's accuracy and effectiveness in predicting GO terms for diverse biological functions.

2.4. Attention- and transformer-based protein function prediction

TALE[29] uses a self-attention-based transformer to extract representative features from protein sequence to improve protein function prediction. It also leverages a zero-shot learning paradigm to jointly embed sequence and hierarchical

function labels into the latent space, allowing a more cohesive representation of the relationships between features and labels. This joint embedding facilitates TALE to generalize well to novel sequences and unseen function by matching similarities among function labels and sequences. Furthermore, TALE introduces a new loss function to address the issue of hierarchical violation. This loss function includes a hierarchical regularization term, which specifically aims to prevent the predicted scores (probabilities) of child GO terms from surpassing those of its ancestors. Additionally, TALE+, a method that ensembles the top three TALE models and a sequence similarity-based protein function prediction method based on DIAMOND [37], was developed to improve the predictions made by TALE.

TEMPROT [38] is another sequence-based protein function prediction method leveraging ProtBERT-BFD[39], a transformer language model pre-trained on the BFD dataset[8, 40, 41]. The pretrained ProtBERT-BFD was first fine-tuned. The fine-tuning process employs a sliding window technique, dividing sequences into 500 chunks to accommodate ProtBERT-BFD's length limitation of 512. After fine-tuning, the backbone of ProtBERT-BFD is used to extract representative features from protein sequences. These features serve as an input for a meta-classifier based on a multi-layer perceptron to predict protein function. Furthermore, TEMPROT+ combining TEMPROT and a sequence-similarity search tool, BLASTp[42], was developed to improve the prediction performance.

SPROF-GO [43] is a sequence-based alignment-free protein function prediction method, which harnesses a pre-trained protein language model for efficient extraction of informative sequence embeddings, while applying self-attention pooling to focus on crucial residues. Its prediction has three main stages. First, the pre-trained protein language model ProtTrans [39] is used to efficiently extract the initial sequence embedding matrix from sequences. The sequence embedding matrix undergoes parallel processing by two multi-layer perceptrons (MLPs) to acquire an attention vector and a more detailed hidden embedding matrix. The hidden embeddings are then normalized to generate an embedding vector, which is used as an input for an MLP to predict the probabilities of GO terms. SPROF-GO also employs a hierarchical learning strategy to guarantee the consistency among predictions. Furthermore, a label diffusion algorithm is integrated in the test phase to exploit the homology information of proteins with related functions.

ATGO[44] harnesses protein language models trained on extensive sequences in an unsupervised fashion to predict protein function. The strategy aims to address the limitations associated with imbalanced annotated functional data. Specifically, ATGO uses the ESM-1b transformer[45] to extract multi-layer feature embeddings from protein sequences. A supervised triplet neural network was trained on these extracted feature embeddings in order to maximize the difference between positive and negative samples. To further enhance ATGO's performance, a composite method, ATGO+ was also introduced. It combines predictions from ATGO and the Sequence Alignment-Based GO Prediction (SAGP).

PANDA2[46] uses a Graph Neural Network (GNN) to model the GO direct acyclic graph (DAG) representing the hierarchical structure of GO terms. It also incorporates features produced by the transformer-based protein language model ESM[45]. PANDA2 has three blocks serving as fundamental building blocks for refining edge, node, and global features. In the first two blocks, it sequentially updates edge features,

node features, and global features by integrating information of all available features in the GNN. Furthermore, it employs a fully connected layer to change the size of ESM features to the number of classes being considered. Then, it merges node features, the output generated by fully connected layer, DIAMOND scores, and priority scores. This comprehensive combination of information is used as input for the third GNN block. The node features of the third GNN block are used by a sigmoid function to predict the probability of each class (GO term). PANDA2 demonstrates the effectiveness of using a GNN architecture for modeling the GO DAG topology and annotating protein functions.

3. Structure-based protein function prediction

The sequence-based function prediction approach has been more common in protein function prediction than the approaches of using other inputs due to the universal availability of protein sequence, even though other data such as protein structure can provide additional complementary information to improve protein function prediction. With the recent development of high-accuracy protein structure prediction tools such as AlphaFold2[8, 16], protein structures have become generally available and started to be used more and more in protein function prediction. Most structure-based prediction methods use various Graph Neural Networks(GNN) such as graph convolutional network (GCN) and Graph Attention Network(GAT) to represent and process protein structures.

DeepFRI[47] relies on a Graph Convolutional Network (GCN)[48] to integrate protein structures and sequence features extracted from a language model to predict protein function. DeepFRI utilizes known protein structures available in the PDB or homology-based structural models built by SWISS-MODEL[49] as structural input. It uses a language model comprised of a long short-term memory(LSTM) network trained in a self-supervised learning manner to extract residue-level features from protein sequences, followed by the GCN layers merging the residue-level features with the graph built from the contact maps calculated from the input protein structure to generate protein-level feature representations. The protein-level features are used to predict GO terms in each of three function categories: Cellular Component, Biological Process, and Molecular Function as well as the Enzyme Commission (EC) numbers, respectively. DeepFRI also employs gradient-weighted Class Activation maps (grad-CAMs) to elevate the representation resolution from protein-level to the region-level, which allows the detection of function-specific structural sites, facilitating the identification of crucial residues correlated with specific functions.

Different from the GCN used by DeepFRI, **GAT-GO**[50] uses a Graph Attention Network (GAT) to integrate both predicted protein structural information and protein sequence embeddings for accurate protein function prediction. The method uses RaptorX[51, 52] to predict protein structural information (i.e., protein contact map) and ESM-1b to generate sequence embeddings. It first uses a one-dimensional CNN to take both sequential features and residue-level sequence embeddings to create per-residue feature embeddings. Then, the CNN-generated embeddings combined with a RaptorX-predicted contact map are fed into GAT which produces an intermediate embedding that captures both sequential and structural information. The representation constructed by GAT

passes through a dense classifier to predict the probability of protein function terms.

Different from DeepFRI and GAT-GO using earlier protein structure prediction methods to generate structural input, **TransFun**[30] uses AlphaFold-predicted protein structures as input. It employs a transformer-based protein language model and rotation- and translation-equivariant graph neural networks (EGNNs) [53] to distill information from both protein sequences and structures to predict protein functions. Its prediction process has three main stages: 1) building a protein graph from a predicted structure, 2) generating the embeddings from a protein sequence, and 3) using an EGNN model to predict protein functions. In the first stage, protein graphs are generated from protein structures collected from AlphaFoldDB[8, 16] using a K-nearest neighbor(KNN) approach based on the distance between carbon-alpha atoms in a protein structure. In the second stage, per-residue and per-sequence embeddings for proteins are generated from protein sequences by the ESM-1b[45] pre-trained language transformer model. In the final stage, both the per-residue and per-sequence features are combined by the EGNNs to predict protein function.

Struct2GO[54] is also a structure-based method that combines sequence features with structural features obtained from AlphaFold2-predicted structures. It extracts a two-dimensional (2D) protein contact map for an input protein from the three-dimensional (3D) protein structure according to a distance threshold of 10Å between carbon-alpha atoms. Additionally, Node2vec[55] algorithm is employed to generate residue-level features for the protein. The contact map serves as the adjacency matrix of the input graph, which are combined with the node features, i.e., the residue-level features, to generate a graph representation of the protein. The representation is used by a Graph Convolution Neural (GCN) network to generate hidden structural features. The feature generation is enhanced with a self-attention mechanism and the integration of sum-pooling and max-pooling techniques. Additional sequence features are also extracted using the SeqVec[56]. Finally, the sequence features are fused with the structural features as input for a final classifier to make function prediction.

4. Interaction-based protein function prediction

Due to the fact that proteins rarely function in isolation, protein-protein interaction information can be used to enhance protein function prediction. It is particularly useful for predicting GO terms describing biological processes that involve multiple proteins cooperating together. Protein function prediction methods relying on protein-protein interactions primarily focus on genome-scale interaction networks, aggregating data from various sources to gain insights into the functional organization of proteins. Some of these methods emphasize the integration of heterogeneous information from diverse interaction networks. A straightforward approach for data integration is to process each network separately and then combine the features generated from each of them. However, this approach often encounters some challenges like increased dimensionality, information loss, and noise accumulation from high-throughput experiments. In this section, we discuss the diverse approaches of integrating multiple heterogeneous networks to predict protein function.

Mashup[57] is an integrative framework designed to extract high-quality and compact topological feature representations from one or more interaction networks constructed from heterogeneous data types. Although Mashup does not inherently use a deep learning technique, it provides a method for extracting features from multiple heterogeneous networks, which are readily used by several interaction-based deep learning methods [58, 59]. The method consists of three main stages: a diffusion stage, an embedding stage, and a learning stage. The diffusion stage involves applying a localized network diffusion technique, specifically Random Walks with Restart (RWR), to each individual network to obtain a matrix representation capturing the interactions between nodes denoting proteins. This captures information about topological structure and connectivity of nodes in each network. Next, the embedding phase focuses on obtaining low-dimensional feature vectors that represent the topology of each node, which is achieved by minimizing the difference between observed diffusion states and parameterized multinomial logistic distributions across all networks. Finally, the learned representations are used as input features for various downstream tasks including protein function prediction.

Following a similar approach as Mashup, **deepNF**[58] integrates diverse heterogeneous protein interaction networks using deep learning techniques. The process begins with the Random Walk with Restart (RWR) algorithm to obtain high-quality vector representations of proteins in each network, capturing their structural information. A Positive Pointwise Mutual Information (PPMI) function is then applied for normalization, and this process is iterated for each network. The subsequent stage focuses on creating a comprehensive representation by integrating the multiple PPMI instances. To achieve this, deepNF employs a Multimodal Data Autoencoder (MDA) network to encode diverse PPMI instances into a representative matrix and reconstruct it through a decoder. The encoder produces low-dimensional non-linear embeddings for each network, and these representations are concatenated. A common feature representation is computed using multiple non-linear functions. In the decoding phase, the process is reversed to compute larger common representations from individual ones, followed by the reconstruction of PPMI matrices for each network. The final step predicts protein functions based on the comprehensive representations obtained in the bottleneck layer of the autoencoder network.

Similar to Mashup and deepNF, **MELISSA**[59] predicts functions from multiple protein-protein interaction networks. However, the integration of known functional labels during the embedding process sets MELISSA apart from the aforementioned methods. Its prediction unfolds in five key steps: Biclustering, Graph Augmentation, Diffusion, Embedding, and Learning. In the initial stage, MELISSA employs a biclustering algorithm to simultaneously cluster proteins and functional labels. This results in biclusters where proteins within clusters share similar functional labels, and functional labels are rarely shared across clusters. In the following step, the protein-protein interaction graphs undergo augmentation by introducing auxiliary nodes, each representing a distinct cluster. Nodes in the graph are then connected to their corresponding auxiliary nodes using must-link constraints (positive weighted edges). Additionally, pairwise cannot-link constraints (edges with negative weights) are introduced between the auxiliary nodes. This augmentation transforms the graphs into signed graphs, where auxiliary nodes encode functional information. Nodes within the same cluster are

drawn closer, while nodes in different clusters are pushed apart. Following the augmentation stage, diffusion state matrices are generated for each augmented graph using a generalization of the method applied in Mashup, by considering the signed nature of the edges. In the final step, MELISSA follows Mashup’s approach to generate embeddings for each node. These embeddings can be effectively utilized by existing function prediction methods to predict function terms.

NetQuilt[60] is a method that integrates protein sequence and protein-protein interaction (PPI) information from multiple species. The approach computes similarity scores between proteins across species using a recurrence equation derived from the IsoRank method of multi-species network alignment [61]. A large symmetric similarity matrix is constructed, where IsoRank similarity matrices of all species with themselves are placed along the diagonal, resulting in a block-diagonal matrix. Interspecies protein similarity matrices are placed on the off-diagonal. The matrix then contains the information from all the individual protein interaction networks as well as the links between them.

The matrix constructed, along with sequence-similarity information, is used as input for a maxout neural network to predict protein function.

DeepGO[62] introduces an approach to predict protein function based on protein sequences and known interactions. It integrates features derived from sequences and protein-protein interaction (PPI) networks across various species in the STRING database. The combined sequence and PPI network features undergo processing in a fully connected layer, and the resultant output feeds into hierarchically structured neural networks to make function prediction.

STRING2GO[63] employs a deep maxout neural network to acquire functional representations by simultaneously encoding both protein-protein interactions and functional annotation information. It uses two methods to generate network embedding representations, (1) a network embedding generation process similar to the one in mashup and (2) node2vec of generating embeddings from the STRING network. After the generation of embeddings, Deep Maxout Neural Networks (DMNNs) is used to simultaneously learn and encode representation information from both the protein-protein interaction network and protein functional annotations. The functional representations are extracted from the outputs of the third hidden layer of DMNNs, which is used by a Support Vector Machine (SVM) to predict the probability of GO terms.

5. Integrative protein function prediction

In this section, we will delve into the methods of integrating multiple sources of information to predict protein function.

DeepGraphGO[64] aims to tackle the limitation of protein interaction-based methods that did not include sequence information. It introduced a multi-species strategy to incorporate the data of all species to train a single model. This approach significantly augments the number of training samples, surpassing the capabilities of existing network-based methods using less data at the time. Binary input protein features are generated through InterProScan, wherein each element indicates the presence or absence of a protein domain, family, or motif. These binary features are combined with protein network graphs, where proteins serve as the nodes and protein-protein interactions form the edges for functional protein annotation. DeepGraphGO prediction

comprises three primary steps. First, a fully connected layer is employed to convert the binary features into a non-binary vector with reduced dimensions, serving as the initial feature representation. Next, updating the representation vector of each node and incorporating new information from network interactions is achieved through a graph convolutional neural network. Finally, a fully connected layer is utilized to predict probabilities of GO terms.

Graph2GO[65] is a multi-modal graph-based representation learning model that integrates heterogeneous information. This model incorporates multiple types of protein interaction networks derived from sequence similarity and protein-protein interaction, along with protein features such as amino acid sequence, subcellular location, and protein domains. The Graph2GO pipeline is composed of two Variational Graph Auto-Encoder (VGAE)[66] models for the protein-protein interaction network and sequence similarity network (SSN). These VGAE models extract representative embeddings, which are subsequently used as input to a final fully-connected deep neural network (DNN) classifier for the prediction of protein functions.

Three version of NetGO methods, **NetGO**, **NetGO2**, and **NetGO3** are related to an early integrative method - **GOlabeler**[67], which encompasses five distinct components: Naive prediction (GO term frequency), BLAST-KNN (k-nearest neighbor using BLAST results), LR-3mer (Logistic regression of the frequency of amino acid trigrams), LR-InterPro (Logistic regression of InterPro features utilizing rich domain, family, and motif information), and LR-ProFET (Logistic regression of ProFET features). The outputs of these components are combined through learning to rank (LTR) to predict protein function. **NetGO**[68] introduces a novel component, Net-KNN, incorporating network information into the system. **NetGO2**[69] further enhances the system by incorporating two additional components, LR-Text and Seq-RNN, while excluding the LR-ProFET component. For LR-Text, corresponding text data about proteins is extracted from PubMed, forming a document that is represented using sparse TF-IDF (term frequency-inverse document frequency) and dense semantic representations generated by Doc2Vec[70]. Logistic regression is trained with these text-based features. Meanwhile, Seq-RNN is employed to extract deep representations of protein sequences, using a Bi-directional Long Short-Term Memory (BiLSTM), followed by a fully connected layer to predict functions. **NetGO3**[71] modifies the architecture by replacing the Seq-RNN component with LR-ESM. LR-ESM generates embeddings for each protein using ESM-1b[45].

SDN2GO[72] employs an integrated deep learning model combining protein sequence, protein domains, and protein-protein interaction networks for protein function prediction. The model has four parts, a sequence sub-model, a domain sub-model, a PPI-net sub-model, and a weighted classifier. The sequence sub-model extracts features from sequence input, which is represented as two-dimensional 3-grams-vector-matrix. The model uses one-dimensional CNNs to extract in-depth high-dimensional features. The PPI-net sub-model utilizes three-layer trapezoidal neural networks to generate the features of PPI Network input. The domain sub-model uses the sorted protein domain information as an input for a sparse layer to generate intermediate features. The output of the Sparse layer represented as two-dimensional matrix enters one-dimensional CNNs to extract features. The output features represented as

vectors with same dimensions generated by all the three sub-models are combined as input for the weighted classifier to predict functions of protein.

PFP-GO[73] also integrates protein sequence, protein domain, and PPI network information for protein function prediction. It first uses the information separately to rank each individual GO term, and the ranking determines which GO terms are associated with the target proteins. In this method, mapping data from one source to another becomes crucial as three complementary information sources are utilized. It makes predictions in four steps. Firstly, a PPI network for target proteins is obtained. Secondly, only the level-2 neighborhood graph for each target protein is taken into account, eliminating other non-essential proteins. Thirdly, after acquiring refined PPI for each target protein, GO terms are assigned to the target protein and its neighbors using the sequence-based, domain-based, and interaction neighbor-based approaches. Lastly, GO terms are ranked based on a function enrichment score, and a consensus score is applied to select GO terms for each target protein.

Like PFP-GO, **MultiPredGO**[74] predicts protein functions by combining protein sequence, protein structure, and PPI network information. Two individual deep learning models are used for feature extraction from sequence and structure, and a pre-trained knowledge graph embedding method is used for PPI network. The sequence is first transformed into a trigram and then processed by an embedding layer. Then, the embedding output passes through one-dimensional convolutional layer for feature extraction. For the structure, a 3D structure is retrieved from Protein Data Bank (PDB) if available, and converted into four distinct 3D voxel representation. Then, an off-shelf residual network, ResNet-50[75], is employed to extract features from the structure. Lastly, extracted features from sequence and structure are combined with PPI network information to obtain the final features, which are processed by a neuro-symbolic hierarchical classifier to make function prediction.

Finally, **DeepGATGO**[76] is an integrative function prediction method leveraging a graph attention learning network(GATs) and a contrastive learning[77, 78] approach to combine protein sequence information and structural and semantic information of Gene Ontology (GO) terms to predict protein functions. It utilizes ESM-1b[45] pre-trained language model to extract feature embeddings from protein sequences. The structural information of GO terms is extracted using GAT network. The semantic information of GO terms is generated through contrastive learning from embeddings created using their names and textual descriptions by the BioBert[79] pre-trained natural language processing model. The extracted semantic features and structural features of GO terms are concatenated. The resulting concatenation output is then multiplied with the protein sequence features. The concatenated features are used by a classification layer with the triplet loss and binary cross-entropy loss to predict the functions of proteins.

6. Few-shot learning-based protein function prediction

One significant challenge in protein function prediction is to predict GO terms that are associated with few proteins because they are severely underrepresented or not present in the training data. For instance, more than 20,000 GO terms have less than 100 annotated proteins possessing them as function.

One way to tackle this problem is to use semantic information of GO terms. Given the scarcity of labeled examples for rare GO terms, semantic information is harnessed to establish meaningful relationships between rare GO terms and common GO terms. Examples of semantic information include leveraging the hierarchical relationships within the GO graph and utilizing GO textual descriptions. Another way is to apply embedding functions to associate features with labels, projecting both feature and label embeddings into a common space and aligning similar GO terms nearby. **TALE**[29] jointly embeds sequence and hierarchical function labels into a latent space, allowing it to generalize to novel/rare terms. Tale focuses on terms that have at least one protein annotation and simultaneously embeds protein sequences and hierarchical function labels using the attention mechanism.

ProTranslator[31] transfers function annotations with similar textual descriptions to annotate a novel function. Leveraging textual descriptions, ProTranslator embeds Gene Ontology (GO) functions using their textual descriptions. The embedding is performed using PubMedBert[80], a language model pre-trained on PubMed abstracts and full-text articles. Proteins are embedded to generate three widely-used features: sequence features, textual description features, and PPI-network features. Similar to deepGOPlus, the sequence features are extracted using convolutional neural networks (CNN) with multiple one-dimensional convolution kernels. Textual descriptions are obtained from GeneCards [21]. The PPI-network features are obtained from pre-trained Mashup representations calculated from protein-protein interaction networks. Ultimately, GO terms and proteins are projected into the same low-dimensional space using a bilinear layer.

DeepGOZero[81] improves predictions for rare GO classes with limited or zero annotations using a model-theoretic approach (ELEmbeddings [82]) to learn ontology embeddings. The ELEmbeddings represent classes as n-balls and relations as vectors to embed ontology semantics into a geometric model. It also uses Interpro domain annotations to generate an embedding of size 1024 for each protein. The protein embeddings and ontology embeddings are combined to predict GO terms.

7. Data Sources, Critical Assessment of Protein Function Annotation (CAFA), and Evaluation Metrics

7.1. Data Sources

Curating high-quality training and test datasets is a key to develop accurate deep learning methods for protein function prediction. Protein sequences and function labels are often sourced from the UniProt Knowledgebase(UniProtKB) [20]. UniProtKB consists of two sections: UniProtKB/Swiss-Prot(reviewed, manually annotated proteins) and UniProtKB /TrEMBL (unreviewed, automatically annotated proteins). The former contains protein sequences and function labels that have been carefully, manually-annotated, while the latter includes computationally analyzed records awaiting full manual annotation. To obtain high-quality labels, the proteins in UniProtKB/Swiss-Prot are usually used to create training and test datasets.

The structure for a protein can be directly predicted by protein structure prediction tools such as AlphaFold or collected from PDB[17] and AlphaFoldDB[16] if available. PPI

networks are usually retrieved from the STRING database integrating huge amounts of experimentally determined and predicted protein-protein interactions. InterPro is a valuable source to obtain the family and function motif/site annotations for proteins and domains, which can be used as input features for protein function prediction. InterPro integrates the data from 13 member databases, forming the InterPro consortium, including CATH[83, 84], CDD[85], HAMAP[86], MobiDB Lite[87], Panther[88], Pfam[89], PIRSF[90], PRINTS[91], Prosite[92], SFLD[93], SMART[94], SUPERFAMILY[95, 96] and NCBIfam. All the features for a protein in Interpro can be obtained using the interproscan (a tool to scan sequences against all InterPro’s member databases) or downloaded from the InterPro website. Finally, protein textual descriptions can be gathered from UniProtKB and GeneCards.

7.2. Critical Assessment of Function Annotation (CAFA)

Objectively and rigorously assessing the performance of different protein function prediction methods is important to advance the field. The Critical Assessment of Function Annotation(CAFA)[11, 12], a global, community-wide experiment held every few years to blindly assess protein function prediction methods. It uses proteins whose function annotations are not available as targets for participating methods to predict their function. The prediction results are then evaluated when the true function annotations of the targets become available. Several CAFA experiments have been held, including the inaugural challenge (CAFA1) taking place in 2010-2011 and the most recent challenge, CAFA5, held in 2023. According to the first four rounds of CAFA experiments (CAFA1-4), the performance of protein function prediction has gradually progressed over years. The results of CAFA5 remain to be seen.

7.3. Evaluation metrics

Evaluating protein function prediction using multiple complementary metrics is important to assess the strength and weakness of function prediction methods. A list of commonly used metrics for evaluating GO term predictions including F_{\max} , AUPR, AUC, MCC, and S_{\min} are briefly discussed below.

F_{\max} is one of the main metrics used in the field as well by CAFA[12, 11]. It is the maximum F-measure score(F1 score: the geometric mean of precision and recall) among all the F1 scores calculated for all the prediction decision thresholds(t), where the precision (Pr) and recall (Rc) for a decision threshold (or cut-off value) t is calculated as follows:

$$Pr(t) = \frac{TP(t)}{TP(t) + FP(t)}$$

where TP is the number of True Positives and FP is the number of False Positives.

$$Rc(t) = \frac{TP(t)}{TP(t) + FN(t)}$$

where TP is the number of True Positives and FN is the number of False Negatives.

The F1 score for a decision threshold t is then computed as follows:

$$F1(t) = 2 \times \frac{Pr(t) \times Rc(t)}{Pr(t) + Rc(t)}$$

Finally, F_{\max} is calculated as the maximum F1 score over all decision thresholds:

$$F_{\max} = \max_t (F1(t))$$

AUPR stands for Area Under the Precision-Recall curve, which is also a commonly used evaluation metric. Similarly, AUC measuring the area under the Receiver Operating Characteristic (ROC) curve is often used. A ROC curve is a plot of the True Positive Rate (TPR) against the False Positive Rate (FPR) across different cut-off values t . TPR and FPR for a cut-off value t is defined as follows:

$$\text{TPR}(t) = \frac{\text{TP}(t)}{\text{TP}(t) + \text{FN}(t)}$$

$$\text{FPR}(t) = \frac{\text{FP}(t)}{\text{FP}(t) + \text{TN}(t)}$$

The Matthews Correlation Coefficient (MCC) is a metric that is particularly useful when test datasets are significantly imbalanced. It is calculated with the following formula:

$$\text{MCC} = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

S_{\min} is another popular evaluation metric used in CAFA[12, 11]. It measures the minimum semantic distance between GO terms based on remaining uncertainty (ru) and misinformation (mi). The remaining uncertainty of the true annotation of protein represents the information that has not been provided or accounted for by the predicted annotation. The misinformation represents a metric that measures the level of misleading information linked to a predicted annotation. The S_{\min} is expressed as follows:

$$S_{\min} = \min_t \left(\sqrt{ru(t)^2 + mi(t)^2} \right)$$

8. Challenges and Future Direction

As discussed in the previous sections, substantial advances in developing deep learning methods for protein function prediction have been made by the community in the last several years. However, the accuracy of protein function still has not reached the high-accuracy level of protein structure prediction that has made it an indispensable tool for biomedical research. There are at least three major challenges in protein function prediction that need to be addressed in order to substantially improve its accuracy.

The first major challenge is to develop highly sophisticated deep learning and AI methods to synergistically integrate multiple modalities of input data (e.g., protein sequence, protein structure, protein interaction, protein/domain family information, and biological textual description) to improve protein function. Most existing integrative methods simply extract features from each data modality and then concatenate them without letting modalities systematically interact with each other in the feature extraction process. The techniques used by the large language models (LLMs) such as ChatGPT-4 and Gemini[97] to integrate multiple modality data such as text, image, video, and voice through seamless cross-modality communication may be transferred to the protein function prediction field to integrate multiple modalities of protein data. And it is time to develop multi-modal large language models for proteins as multi-modality protein data such as sequences and structures are ubiquitously available nowadays.

The second major challenge is how to more effectively leverage the evolutionary information hidden in the hundreds of millions of protein sequences better to improve protein function prediction. A promising direction is to develop more

sophisticated large language models for protein sequences (LLMP) that can be directly fine tuned or promoted to predict protein function. The current application of LLMP such as ESM-1b is still in the early stage and at a shallow level because the pretrained LLMP are mostly used to generate features from sequences as input for protein function prediction. One way to deepen the application of LLMP in protein function prediction is to directly fine tune the weights of the pretrained LLMP component in the protein function prediction system during the training of the system. Another way is to add function prediction into the designing and training of LLMPs in the first place so that they are intrinsically built for protein function prediction. For instance, a LLMP can be designed to predict masked or next amino acids through self-supervised learning as well as function terms through supervised learning. The LLMP can be mainly trained on millions of unlabeled protein sequences to predict masked or next amino acids and auxilinarly trained to predict function terms of thousands of proteins with function labels at the same time as how a large language model for natural language processing (NLP) was trained to predict next (masked) tokens and classify sentences simultaneously [98].

The third major challenge is to improve the prediction accuracy for rare GO terms with low frequency in protein function annotations or novel GO terms that never occur before. Some rare GO terms are highly specific GO terms that occur at the bottom level of the gene ontology graph, which are important for protein function annotation but very hard to predict. As demonstrated by some zero- or few-shot prediction methods such as TALE[29] and ProTranslator[31], zero-shot or few-shot learning methods [99] used in NLP, computer vision and image processing may be transferred to the field of protein function prediction. Particularly, we envision that the prompt engineering and in-context learning [100] used with large language models (LLMs) for NLP can also be used with LLMPs to predict rare or novel GO terms, provided that LLMPs fine-tuned for protein function prediction, akin to LLMs for NLP, are developed in the field. Therefore, a user can use one or a few rare GO terms as examples as prompts to guide the pretrained LLMPs to predict rare or novel GO terms in any context as one uses prompts to instruct ChatGPT to learn new concepts or skills.

In summary, we envision that developing next-generation sophisticated LLMPs that can handle multiple modalities of protein data, be fine tuned directly by function labels, or be customized by prompt-based in-context learning for protein function prediction may be a promising avenue for tackling some major challenges in protein function prediction, such as multi-modality data integration, extracting evolutionary information from millions of sequences, and predicting rare/novel GO terms, to push the performance of protein function prediction to the next level.

9. Competing interests

No competing interest is declared.

10. Author contributions statement

J.C. conceived the review project and the future development directions. F.B. and A.L. collected the data. F.B., A.L., and J.C. wrote the manuscript.

11. Acknowledgments

This work is supported in part by a grant from the National Science Foundation (NSF grant #: DBI2308699).

References

1. A. LaPelusa and R. Kaushik. Physiology, proteins. *StatPearls*, Nov 14 2022. In: StatPearls [Internet]. Treasure Island (FL): StatPearls Publishing; 2023 Jan–.
2. Nabin Giri and Jianlin Cheng. De novo atomic protein structure modeling for cryo-em density maps using 3d transformer and hidden markov model. *bioRxiv*, 2024.
3. Simon C. Bull and Andrew J. Doig. Properties of protein drug target classes. *PLoS One*, 10(3):e0117955, Mar 30 2015.
4. Ronaldo Santos, Oleg Ursu, Anna Gaulton, A. Patrícia Bento, Ricardo S. Donadi, Cristian G. Bologna, Anna Karlsson, Bissan Al-Lazikani, Anne Hersey, Tudor I. Oprea, and John P. Overington. A comprehensive map of molecular drug targets. *Nat Rev Drug Discov*, 16(1):19–34, Jan 2017. Epub 2016 Dec 2.
5. Ashwin Dhakal, Cole McKay, John J Tanner, and Jianlin Cheng. Artificial intelligence in the prediction of protein–ligand interactions: recent advances and future directions. *Briefings in Bioinformatics*, 23(1):bbab476, 11 2021.
6. Andrew W Senior, Richard Evans, John Jumper, James Kirkpatrick, Laurent Sifre, Tim Green, Chongli Qin, Augustin Židek, Alexander WR Nelson, Alex Bridgland, et al. Protein structure prediction using multiple deep neural networks in the 13th critical assessment of protein structure prediction (casp13). *Proteins: structure, function, and bioinformatics*, 87(12):1141–1148, 2019.
7. Jie Hou, Tianqi Wu, Renzhi Cao, and Jianlin Cheng. Protein tertiary structure modeling driven by deep learning and contact distance prediction in casp13. *Proteins: Structure, Function, and Bioinformatics*, 87(12):1165–1178, 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 Dauparas, 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. Predrag Radivojac, Wyatt T Clark, Tal Ronnen Oron, Alexandra M Schnoes, Tobias Wittkop, Artem Sokolov, Kiley Graim, Christopher Funk, Karin Verspoor, Asa Ben-Hur, et al. A large-scale evaluation of computational protein function prediction. *Nature methods*, 10(3):221–227, 2013.
11. Yuxiang Jiang, Tal Ronnen Oron, Wyatt T Clark, Asma R Bankapur, Daniel D’Andrea, Rosalba Lepore, Christopher S Funk, Indika Kahanda, Karin M Verspoor, Asa Ben-Hur, et al. An expanded evaluation of protein function prediction methods shows an improvement in accuracy. *Genome biology*, 17(1):1–19, 2016.
12. Naihui Zhou, Yuxiang Jiang, Timothy R Bergquist, Alexandra J Lee, Balint Z Kacsóh, Alex W Crocker, Kimberley A Lewis, George Georgiou, Huy N Nguyen, Md Nafiz Hamid, et al. The cafa challenge reports improved protein function prediction and new functional annotations for hundreds of genes through experimental screens. *Genome biology*, 20(1):1–23, 2019.
13. Gene Ontology Consortium. The gene ontology (go) database and informatics resource. *Nucleic acids research*, 32(suppl.1):D258–D261, 2004.
14. 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.
15. Farhad Pourpanah, Moloud Abdar, Yuxuan Luo, Xinlei Zhou, Ran Wang, Chee Peng Lim, Xi-Zhao Wang, and QM Jonathan Wu. A review of generalized zero-shot learning methods. *IEEE transactions on pattern analysis and machine intelligence*, 2022.
16. 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.
17. Helen M Berman, John Westbrook, Zukang Feng, Gary Gilliland, Talapady N Bhat, Helge Weissig, Ilya N Shindyalov, and Philip E Bourne. The protein data bank. *Nucleic acids research*, 28(1):235–242, 2000.
18. Damian Szklarczyk, Andrea Franceschini, Stefan Wyder, Kristoffer Forslund, Davide Heller, Jaime Huerta-Cepas, Milan Simonovic, Alexander Roth, Alberto Santos, Kalliopi P Tsafou, et al. String v10: protein–protein interaction networks, integrated over the tree of life. *Nucleic acids research*, 43(D1):D447–D452, 2015.
19. Sarah Hunter, Rolf Apweiler, Teresa K Attwood, Amos Bairoch, Alex Bateman, David Binns, Peer Bork, Ujjwal Das, Louise Daugherty, Lauranne Duquenne, et al. Interpro: the integrative protein signature database. *Nucleic acids research*, 37(suppl.1):D211–D215, 2009.
20. The UniProt Consortium. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids Research*, 51(D1):D523–D531, 11 2022.
21. Gil Stelzer, Naomi Rosen, Inbar Plaschkes, Shahar Zimmerman, Michal Twik, Simon Fishilevich, Tsippi Iny Stein, Ron Nudel, Iris Lieder, Yaron Mazor, et al. The genecards suite: from gene data mining to disease genome sequence analyses. *Current protocols in bioinformatics*, 54(1):1–30, 2016.
22. Keiron O’Shea and Ryan Nash. An introduction to convolutional neural networks. *arXiv preprint arXiv:1511.08458*, 2015.
23. Larry R Medsker and LC Jain. Recurrent neural networks. *Design and Applications*, 5(64-67):2, 2001.
24. Mike Schuster and Kuldeep K Paliwal. Bidirectional recurrent neural networks. *IEEE transactions on Signal Processing*, 45(11):2673–2681, 1997.
25. Vivienne Sze, Yu-Hsin Chen, Tien-Ju Yang, and Joel S Emer. Efficient processing of deep neural networks: A tutorial and survey. *Proceedings of the IEEE*, 105(12):2295–2329, 2017.

26. Sajid Mahmud, Elham Soltanikazemi, Frimpong Boadu, Ashwin Dhakal, and Jianlin Cheng. Deep learning prediction of severe health risks for pediatric covid-19 patients with a large feature set in 2021 barda data challenge. *ArXiv*, 2022.
27. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
28. Jianpeng Cheng, Li Dong, and Mirella Lapata. Long short-term memory-networks for machine reading. *arXiv preprint arXiv:1601.06733*, 2016.
29. Yue Cao and Yang Shen. TALE: Transformer-based protein function Annotation with joint sequence-Label Embedding. *Bioinformatics*, 37(18):2825–2833, September 2021.
30. Frimpong Boadu, Hongyuan Cao, and Jianlin Cheng. Combining protein sequences and structures with transformers and equivariant graph neural networks to predict protein function. *Bioinformatics*, 39(Supplement_1):i318–i325, June 2023.
31. Hanwen Xu and Sheng Wang. Protranslator: zero-shot protein function prediction using textual description. In *International Conference on Research in Computational Molecular Biology*, pages 279–294. Springer, 2022.
32. R. Cao, C. Freitas, L. Chan, M. Sun, H. Jiang, and Z. Chen. Prolango: Protein function prediction using neural machine translation based on a recurrent neural network. *Molecules*, 22(10):1732, 2017.
33. Chang Wook Ko, Juyeon Huh, and Jae Woo Park. Deep learning program to predict protein functions based on sequence information. *MethodsX*, 9:101622, Jan 15 2022.
34. Maxat Kulmanov and Robert Hoehndorf. Deepgoplus: improved protein function prediction from sequence. *Bioinformatics*, 36(2):422–429, January 2020.
35. Wei Xia, Lei Zheng, Jiahua Fang, Feng Li, Yijun Zhou, Zhichao Zeng, Bo Zhang, Zhihao Li, Hui Li, and Fan Zhu. Pfmuld: a novel strategy enabling multi-class and multi-label protein function annotation by integrating diverse deep learning methods. *Computers in Biology and Medicine*, 145:105465, June 2022.
36. Ahmet Sureyya Rifaioğlu, Tunca Doğan, Maria Jesus Martin, and et al. Deepred: Automated protein function prediction with multi-task feed-forward deep neural networks. *Scientific Reports*, 9:7344, 2019.
37. Benjamin Buchfink, Chao Xie, and Daniel H Huson. Fast and sensitive protein alignment using diamond. *Nature methods*, 12(1):59–60, 2015.
38. Gustavo B. Oliveira, Helio Pedrini, and Zanoni Dias. TEMPROT: Protein Function Annotation Using Transformers Embeddings and Homology Search. *BMC Bioinformatics*, 24(1):242, Jun 8 2023.
39. Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Wang Yu, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, Debsindhu Bhowmik, and Burkhard Rost. Prottrans: Towards cracking the language of life's code through self-supervised deep learning and high performance computing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.
40. Martin Steinegger, Milot Mirdita, and Johannes Söding. Protein-level assembly increases protein sequence recovery from metagenomic samples manifold. *Nature methods*, 16(7):603–606, 2019.
41. Martin Steinegger and Johannes Söding. Clustering huge protein sequence sets in linear time. *Nature communications*, 9(1):2542, 2018.
42. SF Altschul, TL Madden, AA Schäffer, J Zhang, Z Zhang, W Miller, and DJ Lipman. Gapped blast and psi-blast: a new generation of protein database search programs. *Nucleic Acids Res*, 25(17):3389–3402, 1997.
43. Qiang Yuan, Jia Xie, Jia Xie, Hui Zhao, and Yu Yang. Fast and accurate protein function prediction from sequence through pretrained language model and homology-based label diffusion. *Brief Bioinform*, 24(3):bbad117, May 19 2023.
44. Yi-Heng Zhu, Chengxin Zhang, Dong-Jun Yu, and Yang Zhang. Integrating unsupervised language model with triplet neural networks for protein gene ontology prediction. *PLOS Computational Biology*, 18(12):1–26, 12 2022.
45. Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *PNAS*, 2019.
46. Chenguang Zhao, Tong Liu, and Zheng Wang. Panda2: protein function prediction using graph neural networks. *NAR Genomics and Bioinformatics*, 4(1):lqac004, March 2022.
47. V. Gligorijević, P.D. Renfrew, T. Kosciolk, et al. Structure-based protein function prediction using graph convolutional networks. *Nature Communications*, 12(1):3168, 2021.
48. Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.
49. Torsten Schwede, Jurgen Kopp, Nicolas Guex, and Manuel C Peitsch. Swiss-model: an automated protein homology-modeling server. *Nucleic acids research*, 31(13):3381–3385, 2003.
50. B. Lai and J. Xu. Accurate protein function prediction via graph attention networks with predicted structure information. *Briefings in Bioinformatics*, 23(1):bbab502, 2022.
51. Jian Peng and Jinbo Xu. Raptorx: exploiting structure information for protein alignment by statistical inference. *Proteins: Structure, Function, and Bioinformatics*, 79(S10):161–171, 2011.
52. J. Xu, M. Mcpartlon, and J. Li. Improved protein structure prediction by deep learning irrespective of co-evolution information. *Nat Mach Intell*, 3:601–609, July 2021.
53. Victor Garcia Satorras, Emiel Hoogeboom, and Max Welling. E (n) equivariant graph neural networks. In *International conference on machine learning*, pages 9323–9332. PMLR, 2021.
54. Peishun Jiao, Beibei Wang, Xuan Wang, Bo Liu, Yadong Wang, and Junyi Li. Struct2GO: protein function prediction based on graph pooling algorithm and AlphaFold2 structure information. *Bioinformatics*, 39(10):btad637, 10 2023.
55. Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 855–864, 2016.

56. Michael Heinzinger, Ahmed Elnaggar, Yu Wang, Christian Dallago, Dmitrii Nechaev, Florian Matthes, and Burkhard Rost. Modeling aspects of the language of life through transfer-learning protein sequences. *BMC bioinformatics*, 20(1):1–17, 2019.
57. Hyunghoon Cho, Bonnie Berger, and Jian Peng. Compact integration of multi-network topology for functional analysis of genes. *Cell systems*, 3(6):540–548, 2016.
58. Vladimir Gligorijević, Meet Barot, and Richard Bonneau. deepNF: deep network fusion for protein function prediction. *Bioinformatics*, 34(22):3873–3881, 06 2018.
59. Kaiyi Wu, Di Zhou, Donna Slonim, Xiaozhe Hu, and Lenore Cowen. Melissa: Semi-supervised embedding for protein function prediction across multiple networks. *bioRxiv*, 2023.
60. Meet Barot, Vladimir Gligorijević, Kyunghyun Cho, and Richard Bonneau. Netquilt: deep multispecies network-based protein function prediction using homology-informed network similarity. *Bioinformatics*, 37(16):2414–2422, August 2021.
61. Rohit Singh, Jinbo Xu, and Bonnie Berger. Global alignment of multiple protein interaction networks with application to functional orthology detection. *Proceedings of the National Academy of Sciences*, 105(35):12763–12768, 2008.
62. Maxat Kulmanov, Mohammed Asif Khan, and Robert Hoehndorf. Deepgo: predicting protein functions from sequence and interactions using a deep ontology-aware classifier. *Bioinformatics*, 34(4):660–668, February 2018.
63. Cen Wan, Domenico Cozzetto, Rui Fa, and David T Jones. Using deep maxout neural networks to improve the accuracy of function prediction from protein interaction networks. *PLoS one*, 14(7):e0209958, 2019.
64. Ronghui You, Shuwei Yao, Hiroshi Mamitsuka, and Shanfeng Zhu. DeepGraphGO: graph neural network for large-scale, multispecies protein function prediction. *Bioinformatics*, 37(Supplement_1):i262–i271, 07 2021.
65. Kunjie Fan, Yuanfang Guan, and Yan Zhang. Graph2GO: a multi-modal attributed network embedding method for inferring protein functions. *GigaScience*, 9(8):giaa081, 08 2020.
66. Thomas N Kipf and Max Welling. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308*, 2016.
67. Ronghui You, Zihan Zhang, Yi Xiong, Fengzhu Sun, Hiroshi Mamitsuka, and Shanfeng Zhu. GOLabeler: improving sequence-based large-scale protein function prediction by learning to rank. *Bioinformatics*, 34(14):2465–2473, 03 2018.
68. 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, 05 2019.
69. 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.
70. Quoc Le and Tomas Mikolov. Distributed representations of sentences and documents. In *International conference on machine learning*, pages 1188–1196. PMLR, 2014.
71. 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*, 2023.
72. Yudong Cai, Jia Wang, and Lei Deng. Sdn2go: An integrated deep learning model for protein function prediction. *Frontiers in Bioengineering and Biotechnology*, 8:391, April 2020.
73. Krishanu Sengupta, Sovan Saha, Asif K Halder, Pritam Chatterjee, Mita Nasipuri, Sudip Basu, and Dariusz Plewczynski. Pfp-go: Integrating protein sequence, domain and protein-protein interaction information for protein function prediction using ranked go terms. *Frontiers in Genetics*, 13:969915, September 2022.
74. Swagarika Jaharlal Giri, Pratik Dutta, Parth Halani, and Sriparna Saha. Multipredgo: Deep multi-modal protein function prediction by amalgamating protein structure, sequence, and interaction information. *IEEE Journal of Biomedical and Health Informatics*, 25(5):1832–1838, 2021.
75. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
76. Zihao Li, Changkun Jiang, and Jianqiang Li. Deepgatgo: A hierarchical pretraining-based graph-attention model for automatic protein function prediction. *arXiv preprint arXiv:2307.13004*, 2023.
77. Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
78. Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673, 2020.
79. Jinhuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 09 2019.
80. 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, 2020.
81. Maxat Kulmanov and Robert Hoehndorf. DeepGOZero: improving protein function prediction from sequence and zero-shot learning based on ontology axioms. *Bioinformatics*, 38(Supplement_1):i238–i245, 06 2022.
82. Maxat Kulmanov, Wang Liu-Wei, Yuan Yan, and Robert Hoehndorf. El embeddings: Geometric construction of models for the description logic el++. *arXiv preprint arXiv:1902.10499*, 2019.
83. Tony E Lewis, Ian Sillitoe, Natalie Dawson, Su Datt Lam, Tristan Clarke, David Lee, Christine Orengo, and Jonathan Lees. Gene3D: Extensive prediction of globular domains in proteins. *Nucleic Acids Research*, 46(D1):D435–D439, 11 2017.
84. Ian Sillitoe, Nicola Bordin, Natalie Dawson, Vaishali P Waman, Paul Ashford, Harry M Scholes, Camilla S M Pang, Laurel Woodridge, Clemens Rauer, Neeladri Sen, Mahnaz Abbasian, Sean Le Cornu, Su Datt Lam, Karel Berka, Ivana Hutařová Varekova, Radka Svobodova, Jon Lees, and Christine A Orengo. CATH: increased structural

- coverage of functional space. *Nucleic Acids Research*, 49(D1):D266–D273, 11 2020.
85. Jiyao Wang, Farideh Chitsaz, Myra K Derbyshire, Noreen R Gonzales, Marc Gwadz, Shennan Lu, Gabriele H Marchler, James S Song, Narmada Thanki, Roxanne A Yamashita, Mingzhang Yang, Dachuan Zhang, Chanjuan Zheng, Christopher J Lanczycki, and Aron Marchler-Bauer. The conserved domain database in 2023. *Nucleic Acids Research*, 51(D1):D384–D388, 12 2022.
 86. Ivo Pedruzzi, Catherine Rivoire, Andrea H. Auchincloss, Elisabeth Coudert, Guillaume Keller, Edouard de Castro, Delphine Baratin, Béatrice A. Cuhe, Lydie Bougueleret, Sylvain Poux, Nicole Redaschi, Ioannis Xenarios, and Alan Bridge. HAMAP in 2015: updates to the protein family classification and annotation system. *Nucleic Acids Research*, 43(D1):D1064–D1070, 10 2014.
 87. Marco Necci, Damiano Piovesan, Damiano Clementel, Zsuzsanna Dosztányi, and Silvio C E Tosatto. MobiDB-lite 3.0: fast consensus annotation of intrinsic disorder flavors in proteins. *Bioinformatics*, 36(22-23):5533–5534, 12 2020.
 88. Huaiyu Mi, Anushya Muruganujan, John T Casagrande, and Paul D Thomas. Large-scale gene function analysis with the panther classification system. *Nature protocols*, 8(8):1551–1566, 2013.
 89. 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.
 90. Cathy H Wu, Anastasia Nikolskaya, Hongzhan Huang, Lai-Su L Yeh, Darren A Natale, Cholanayakanahalli R Vinayaka, Zhang-Zhi Hu, Raja Mazumder, Sandeep Kumar, Panagiotis Kourtesis, et al. Pirsf: family classification system at the protein information resource. *Nucleic acids research*, 32(suppl.1):D112–D114, 2004.
 91. Teresa K Attwood, Alain Coletta, Gareth Muirhead, Athanasia Pavlopoulou, Peter B Philippou, Ivan Popov, Carlos Roma-Mateo, Athina Theodosiou, and Alex L Mitchell. The prints database: a fine-grained protein sequence annotation and analysis resource—its status in 2012. *Database*, 2012:bas019, 2012.
 92. Christian JA Sigrist, Edouard De Castro, Lorenzo Cerutti, Béatrice A Cuhe, Nicolas Hulo, Alan Bridge, Lydie Bougueleret, and Ioannis Xenarios. New and continuing developments at prosite. *Nucleic acids research*, 41(D1):D344–D347, 2012.
 93. Eyal Akiva, Shoshana Brown, Daniel E. Almonacid, 2nd Barber, Alan E., Ashley F. Custer, Michael A. Hicks, Conrad C. Huang, Florian Lauck, Susan T. Mashiyama, Elaine C. Meng, David Mischel, John H. Morris, Sunil Ojha, Alexandra M. Schnoes, Doug Stryke, Jeffrey M. Yunes, Thomas E. Ferrin, Gemma L. Holliday, and Patricia C. Babbitt. The Structure–Function Linkage Database. *Nucleic Acids Research*, 42(D1):D521–D530, 11 2013.
 94. Jörg Schultz, Richard R Copley, Tobias Doerks, Chris P Ponting, and Peer Bork. Smart: a web-based tool for the study of genetically mobile domains. *Nucleic acids research*, 28(1):231–234, 2000.
 95. Julian Gough, Kevin Karplus, Richard Hughey, and Cyrus Chothia. Assignment of homology to genome sequences using a library of hidden markov models that represent all proteins of known structure. *Journal of molecular biology*, 313(4):903–919, 2001.
 96. D Wilson, R Pethica, Y Zhou, C Talbot, C Vogel, M Madera, C Chothia, and J Gough. Superfamily—comparative genomics, datamining and sophisticated visualisation. *Nucleic Acids Res*, 37(D380–D386):14, 2009.
 97. Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
 98. 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.
 99. Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. Generalizing from a few examples: A survey on few-shot learning. *ACM computing surveys (csur)*, 53(3):1–34, 2020.
 100. Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas C Schmidt. A prompt pattern catalog to enhance prompt engineering with chatgpt. *arXiv preprint arXiv:2302.11382*, 2023.