

Protein Function Prediction Through Multi-view Multi-label Latent Tensor Reconstruction

Robert Ebo Armah-Sekum^{1*}, Sandor Szedmak¹ and Juho Rousu^{1*}

¹Department of Computer Science, Aalto University, Konemiehentie 2,
Espoo, 02150, Finland.

*Corresponding author(s). E-mail(s): robert.armah-sekum@aalto.fi;
juho.rousu@aalto.fi;

Contributing authors: sandor.szedmak@aalto.fi;

Abstract

Background: In last two decades, the use of high-throughput sequencing technologies has accelerated the pace of discovery of proteins. However, due to the time and resource limitations of rigorous experimental functional characterization, the functions of a vast majority of them remain unknown. As a result, computational methods offering accurate, fast and large-scale assignment of functions to new and previously unannotated proteins are sought after. Leveraging the underlying associations between the multiplicity of features that describe proteins could reveal functional insights into the diverse roles of proteins and improve performance on the automatic function prediction task.

Results: We present GO-LTR, a multi-view multi-label prediction model that relies on a high-order tensor approximation of model weights combined with non-linear activation functions. The model is capable of learning high-order relationships between multiple input views representing the proteins and predicting high-dimensional multi-label output consisting of protein functional categories. We demonstrate the competitiveness of our method on various performance measures. Experiments show that GO-LTR learns polynomial combinations between different protein features, resulting in improved performance. Additional investigations establish GO-LTR's practical potential in assigning functions to proteins under diverse challenging scenarios: very low sequence similarity to previously observed sequences, rarely observed and highly specific terms in the gene ontology.

Implementation: The code and data used for training GO-LTR is available at <https://github.com/aalto-ics-kepaco/GO-LTR-prediction>

Keywords: protein function, machine learning, CAFA, gene ontology

1 Introduction

As one of the essential biomolecules in living cells, proteins perform a wide range of important functions including aiding cell division, supporting metabolism and providing immune response [18, 22]. Thus, a proficient knowledge of their functions is of crucial biological relevance, especially in elucidating metabolic pathways, understanding disease mechanisms and developing potent drugs. However, of the hundreds of millions of proteins that have been discovered and sequenced using high-throughput technologies, only a small proportion (< 1%) have been functionally characterized [6]. The huge disparity is mainly due to the time and resource constraints of experimental characterization techniques. As a result, computational methods offering fast, large-scale and accurate assignment of functional annotations are highly sought to bridge this ever-widening gap [6, 10].

Proteins are described by several characteristics ranging from the primary sequence, secondary structure, tertiary structure, chemical properties, to the physical interactions they have with other proteins in the performance of their functions [18, 22]. Consequently, several methods utilizing different protein feature sets have been developed, either in a single view or a multi-view setup [17, 28, 44]. Several approaches exist for integrating multiple input views including early, intermediate and late fusion methods. Current function prediction methods have used separate modules to learn salient features from respective feature sets [16, 21, 42] and merged the per-feature representations using concatenation or ranked the terms predicted by each feature-component method.

To actively advance the course of developing computational techniques for protein function annotation, the Critical Assessment of Functional Annotation (CAFA) challenge was introduced a decade ago [28]. Through this initiative, computational models developed are systematically assessed based on their accuracy in assigning functional annotations to new and previously uncharacterized proteins, on benchmark datasets curated from wet lab experiments. Evaluation is done using robust and standardized metrics developed by the community [15, 17, 44]. In CAFA, the Gene Ontology (GO), the most comprehensive resource for protein function annotations [3], is used to represent the functional categories and to evaluate the machine learning models. The thousands of GO categories give rise to a multi-label prediction problem where several GO categories may be valid for a single protein.

In this study, we address the function prediction problem by modeling the joint interactions between different features using the latent tensor reconstruction approach [35, 41], which can be viewed as an extension of higher-order factorization machines [7]. Building on the factorized parameterization and expressivity of factorization machines [7, 29] in modeling complex interactions between variables, LTR leverages the linear form factorization of tensors [19] and a mini-batch data processing scheme, thereby scaling to large datasets while maintaining constant memory and linear time complexity in the size of the input features as well as in the order, size and rank of the tensor. In summary, the study makes the following contributions:

- We present GO-LTR, a multi-view multi-label prediction model for automatic function annotation, based on the latent tensor reconstruction approach.

- We show that GO-LTR improves performance on the function prediction task, as assessed by multiple evaluation metrics.
- We show that leveraging multiple protein modalities, including the recent foundation models based on large language models, results in enhanced performance.
- We present detailed studies on the prediction performance, in terms of similarity to observed sequences in the training set, depth and frequency of GO classes, as well as the prediction threshold of the models.

2 Methods

2.1 Learning Task

In the protein function prediction task, we are given a dataset $\{\mathcal{D} = (\mathbf{x}_i, \mathbf{y}_i) \mid i = 1, \dots, m\}$. For each data sample, there are n_d feature vectors, called the views, $\mathbf{x}_i^{(d)} \mid \mathbf{x}_i^{(d)} \in \mathbb{R}^{n_{x_d}}, d = 1, \dots, n_d$ of potentially different dimensions n_{x_d} , representing different representations of a protein, for example, sequence embeddings, InterPro fingerprints and protein-protein interaction embeddings. Each data sample is associated with a multi-label target vector, $\mathbf{y}_i \in \{0, 1\}^{n_y}$ denoting the membership of the example \mathbf{x}_i to the different functional categories, which in this work are taken from the Gene Ontology (GO).

In the matrix representation, each sample $\mathbf{x}_i^{(d)}$ is embedded into a row of the input matrix $\mathbf{X}^{(d)}$ belonging to view d and \mathbf{y}_i is a row of the label matrix \mathbf{Y} .

2.2 Latent tensor reconstruction (LTR)

We used the latent tensor reconstruction (LTR) model [35, 41] in our experiments (Table 1). LTR, similarly to higher-order factorization machines [7], is based on a tensor-based approximation of a degree n_d polynomial function:

$$f(\mathbf{x}) = \sum_{j=1}^n w_j x_j + \sum_{j,k=1}^n w_{jk} x_j x_k + \dots + \sum_{j_1, j_2, \dots, j_{n_d}=1}^n w_{j_1, \dots, j_{n_d}} x_{j_1}, \dots, x_{j_{n_d}} \quad (1)$$

As the number of parameters in the model is exponential in the polynomial degree n_d , instead a factorized representation (Table 1c), where each regression coefficient w_{j_1, \dots, j_r} is approximated by a weighted sum of products of factor weights, is used:

$$w_{j_1, \dots, j_r} = \sum_{t=1}^{n_t} \lambda_t p_{j_1, t} \cdots p_{j_r, t} \quad (2)$$

This trick provides an exponential reduction in the number of parameters that need to be estimated with both statistical and computational benefits.

In LTR, the parameter tensor $\mathbf{T} = \sum_{t=1}^{n_t} \lambda_t \otimes_{d=1}^{n_d} \mathbf{p}_t^{(d)}$ collecting all the regression coefficients w_{j_1, \dots, j_r} is represented in factorized form as weighted sum of rank-one tensors (Figure 1e). The factor matrices \mathbf{P} containing the factor weights of individual

Table 1 Computations in LTR model. $[m]$ denotes the set $\{1, \dots, m\}$, m refers to the number of data examples, $\langle \cdot \rangle$ denotes the inner product, and $\|\cdot\|$ represents the norm operator. \otimes denotes the tensor product of vectors and \circ connotes the pointwise multiplication of tensors of the same dimension. We use \mathbf{y} to denote a vector and \mathbf{Y} to represent a matrix.

Process	Computations
(a) Multiview data	Given: a sample $\mathcal{S} = ((\mathbf{x}_i^{(1)}, \dots, \mathbf{x}_i^{(n_d)}), \mathbf{y}_i) \mid i \in [m]$, $\mathbf{x}_i^{(d)} \in \mathbb{R}^{n_{x_d}}, d \in [n_d], \mathbf{y}_i \in \mathbb{R}^{n_y}$ Output: Parameter tensor \mathbf{T}
(b) Polynomial regression	$\min_{\mathbf{T}} \sum_i \ y_i - \langle \mathbf{T}, \otimes_{d=1}^{n_d} \mathbf{x}_i^{(d)} \rangle\ ^2, \text{ scalar-valued case}$
(c) Tensor factorization, first level	$\begin{aligned} \mathbf{T} &= \sum_{t=1}^{n_t} \lambda_t \otimes_{d=1}^{n_d} \mathbf{p}_t^{(d)} \\ \pi(\mathbf{x}) &= \sum_{t=1}^{n_t} \lambda_t \langle \otimes_{d=1}^{n_d} \mathbf{p}_t^{(d)}, \otimes_{d=1}^{n_d} \mathbf{x}^{(d)} \rangle = \sum_{t=1}^{n_t} \lambda_t \prod_{d=1}^{n_d} \langle \mathbf{p}_t^{(d)}, \mathbf{x}^{(d)} \rangle \\ &= \mathbf{1}_{n_t}^T \mathbf{D}_\lambda \circ_{d=1}^{n_d} \mathbf{P}^{(d)} \mathbf{x}^{(d)}, \mathbf{D}_\lambda = \text{diag}(\lambda_1, \dots, \lambda_t) \end{aligned}$
(d) Tensor factorization, second level	$\begin{aligned} \mathbf{P}^{(d)} &= \mathbf{V}^{(d)} \mathbf{U}^{(d)T} \mathbf{D}_{\lambda_U}^{(d)}, \quad \mathbf{P}^{(d)} \in \mathbb{R}^{n_t \times n_{x_d}} \\ \ \mathbf{V}_i^{(d)}\ _2 &= 1, i = [n_t], \ \mathbf{U}_j^{(d)}\ _2 = 1, j \in [n_{x_d}], \\ &\mathbf{D}_{\lambda_U} \text{ diagonal} \\ \pi(\mathbf{x}) &= \mathbf{1}_{n_t}^T \mathbf{D}_\lambda \circ_{d=1}^{n_d} (\mathbf{V}^{(d)} \mathbf{U}^{(d)T} \mathbf{D}_{\lambda_U}^{(d)} \mathbf{x}) \end{aligned}$
(e) Vector output for multi-labels	$\pi(\mathbf{x}) = \mathbf{Q}^T \mathbf{D}_\lambda \circ_{d=1}^{n_d} (\mathbf{V}^{(d)} \mathbf{U}^{(d)T} \mathbf{D}_{\lambda_U}^{(d)} \mathbf{x}), \mathbf{Q} \in \mathbb{R}^{n_y \times n_t}$
(f) Including activation functions, e.g., ReLU	$\pi(\mathbf{x}) = \mathbf{Q}^T \mathbf{D}_\lambda \circ_{d=1}^{n_d} \mathcal{B}(\mathbf{V}^{(d)} \mathcal{A}(\mathbf{U}^{(d)T} \mathbf{D}_{\lambda_U}^{(d)} \mathbf{x}))$
(g) Inference	Given: data: \mathbf{x} , parameters: $\mathbf{Q}, \mathbf{D}_\lambda, (\mathbf{V}^{(d)}, \mathbf{U}^{(d)}, \mathbf{D}_{\lambda_U}^{(d)})$, $d = [n_d]$, Output: $\hat{\mathbf{y}} = \pi(\mathbf{x})$
(h) Optimization objective	$\begin{aligned} &\min_{\mathbf{Q}, \boldsymbol{\lambda}, \mathbf{V}^{(d)}, \mathbf{U}^{(d)}, \boldsymbol{\lambda}_U^{(d)}, d \in [n_d]} \frac{1}{2mn_y} \ \mathbf{Y} - \hat{\mathbf{Y}}\ _F^2 + \frac{C_\lambda}{2n_t} \ \boldsymbol{\lambda}\ _2^2 \\ &\text{s.t. } \ \mathbf{Q}_i\ _2 = 1, i \in [n_t], \ \mathbf{V}_i^{(d)}\ _2 = 1, i \in [n_t], \\ &\quad \ \mathbf{U}_j^{(d)}\ _2 = 1, j \in [n_{x_d}], \mathbf{D}_{\lambda_U} \text{ diagonal} \end{aligned}$

variables representation are further factorized through a singular value decomposition (Table 1d, Figure 1f). This reparameterization has the effect of decoupling the factors representing individual variables and further decreasing the number of parameters to estimate.

LTR is capable of handling several, potentially heterogeneous data sources describing the same phenomenon, in a multi-view learning framework (Table 1a). For example, in

the 3-view case studied in the experiments (Section 3.1), cross-view interactions are modeled by the tensor product between the feature vectors of the views.

In the architecture, we introduce further non-linearity to enhance the representation power of the model by the use of Rectified Linear Unit (ReLU) activation functions, \mathcal{A} and \mathcal{B} (Table 1f), applied on the linear layers. Notably, the use of activation functions generalizes the LTR model beyond polynomial functions, such as represented in Equation 1.

To address the multi-label output problem, the vector $\mathbf{1}_{n_t}$ of (Table 1c) is replaced by the learned matrix \mathbf{Q} in (Table 1e), which projects the vector-valued predictions into the output space. Finally, the optimization problem solves a regularised mean squared error between the ground truth \mathbf{Y} and prediction $\hat{\mathbf{Y}}$ (Table 1h).

2.3 Input Data

In this study, we used the “*go-basic.obo*” ontology file [11] released on 1.1.2023, containing information about 46,739 terms. Additionally, we used the manually reviewed and annotated Swiss-Prot protein sequences in the Universal Protein Knowledgebase (UniprotKB) [5] which contains about half a million sequences (Figure 1a). Following the CAFA rules for datasets curation [28], we present two time-separated datasets. Dataset-1 contains protein sequences annotated from the inception of UniprotKB up to 13.03.2023. Dataset-2 on the other hand is a collection of sequences annotated from the 14.03.2023 up to 24.01.2024. The data was filtered to remove duplicate sequences. Next, we selected sequences having at least one annotation supported by any of the following evidence codes: EXP, IDA, IPI, IMP, IGI, IEP, TAS, IC, HTP, HDA, HMP, HGI and HEP.

2.4 Output data

The Gene Ontology (GO) provides a formal and comprehensive representation of the functions of gene products in living organisms using standardized and unified terminology [3]. Functions are described using three main subontologies representing three ancestral nodes: Molecular Function Ontology (MFO), Cellular Component Ontology (CCO) and Biological Process Ontology (BPO).

Annotations are represented in a hierarchical format using a directed acyclic graph (DAG). Within this concept hierarchy, links between nodes are described using relations such as *is-a*, *part-of*, *negatively-regulates* and *capable-of*. Functional terms are related by the true-path propagation rule [4] — where a protein annotated to a deep-level node in the graph is automatically annotated to all its parent terms including the ancestral node(s). This implies that the set of functions associated with a particular protein forms a consistent sub-graph in the DAG.

In this work, we only considered *is-a* relationships and used the true path rule to propagate experimental annotations up to the root terms in the GO graph. The resulting sub-graph is then represented as a binary multi-label target vector. Functional terms having at least 30, 30 and 60 sequence examples were chosen as the final labels for MFO, CCO and BPO respectively. Table 2 provides a summary statistics of the final datasets for all three subontologies.

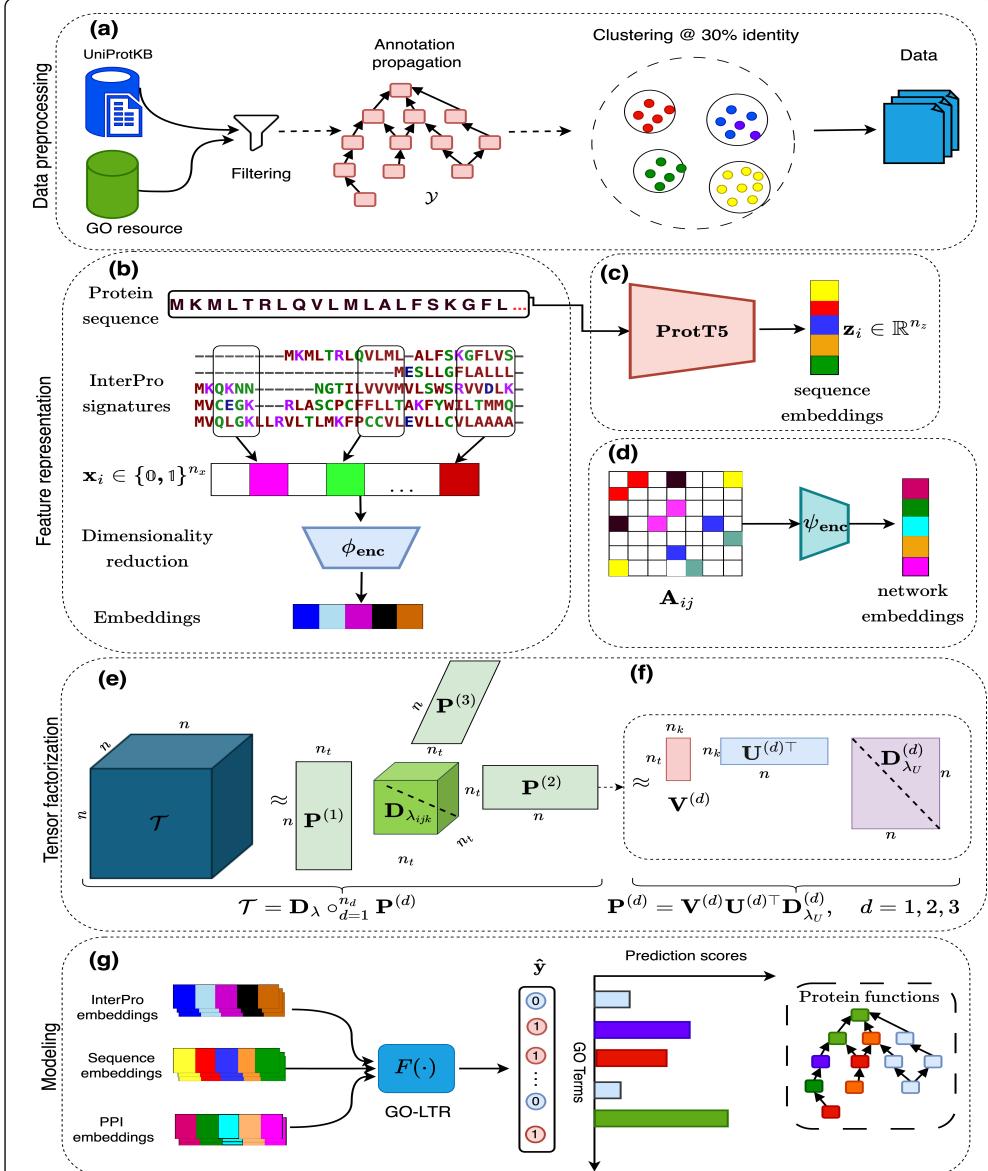


Figure 1 Project workflow: (a) Manually reviewed and annotated sequences are filtered based on term frequency after true-path propagation to ancestor terms. Using MMseqs2, sequences are clustered at 30% percentage identity. (b) Dimensionality reduction is performed on the binary vector of domain and family fingerprints from InterProScan. (c) Sequence embeddings of size 1024 are generated using ProtT5 protein language model. (d) Dimensionality reduction is applied on the rows of the adjacency matrix formed from the PPI network obtained from StringDB. (e) Parameter tensor decomposition in GO-LTR with $D_{\lambda_{ijk}}$ as a diagonal tensor, (f) further singular value decomposition of the parameter matrix $\mathbf{P}^{(d)}$ of each feature (view d). (g) GO-LTR predicts the scores associated with each functional term in the multi-label target vector using multiple features as input. The predicted score for each term is then propagated via the true-path rule to ensure prediction consistency such that each node retains the maximum score during the propagation process.

Table 2 Summary of datasets. Terms do not include the ancestral nodes in the ontology.

Ontology	Terms	Size	
		Dataset-1	Dataset-2
MFO	776	29,650	849
CCO	615	37,025	785
BPO	3,049	39,194	889

2.5 Feature representation

We leveraged three data sources in our experiments — sequence embeddings, interpro fingerprints, and protein-protein interaction data.

Interpro fingerprints [2, 23, 24], encoding information about the motifs, active sites, conserved regions and protein families, were obtained from the Interpro service in the UniprotKB service. This is a binary fingerprint feature describing whether a particular subsequence/domain is present or absent in a sequence. We used an autoencoder composed of four layers interspersed with ReLU activation functions in both the encoder and decoder blocks to reduce the binary feature vector’s dimension from the highly sparse $\approx 14k$ to a dense representation vector of size 1000 (Figure 1b).

We utilized sequence embeddings (1024 dimensions) generated using the ProtT5 [14] language model (Figure 1c). Due to the computationally intensive nature of generating such dense representations, we downloaded the precomputed embeddings made available in the UniprotKB service for all sequences in our dataset. ProtT5 (**P**rotein **T**ext-to-**T**ext **T**ransfer **T**ransformer) is a protein language model developed in [14] based on the transformer architecture [37]. Analogous to Natural Language Processing (NLP) models, ProtT5 considers each input amino acid (AA) sequence as a sentence and its constituent residues as tokens. It is trained in a self-supervised manner: learning to generate the sequence from the low-dimensional intermediate representations from an encoder (see Supplementary Figure A11). The training data for ProtT5 spans over 300 billion amino acids from sequences sourced from large scale databases including Big Fantastic Database (BFD) [12], UniRef50 [34] and UniRef100 [34]. Learning from only the sequence data, the latent representations given by ProtT5 encodes relevant information about proteins including domain, motifs and biophysical features compared to the expensive computation of multiple sequence alignments over large databases.

We also incorporated protein-protein interaction (PPI) network data from the StringDB database [36, 38]. The edges in this network denote the interactions between proteins in the performance of their functions, in a living cell. We create an $N \times N$ adjacency matrix of the PPI graph using all N proteins in our dataset as nodes. Using an autoencoder, we performed dimensionality reduction on each row of the matrix to obtain a 1×1000 dense representation vector (Figure 1d).

2.6 Baseline methods

In addition to commonly used baselines in the automatic function prediction tasks, we also considered machine learning methods that had an open-source implementation that we could train from scratch on our dataset.

BLAST - Basic Local Alignment Search Tool

This method transfers annotations from sequences in the training set to the test set using sequence similarity computed from an optimal alignment between two sequences [1, 25, 32]. Spurious sequence alignments are filtered using an e-value of 0.001. We consider two variants below.

- (i) **BLAST-full** — we transfer all annotations of the training sequence with the highest scoring alignment to a test sample via BLAST, as the multi-label prediction for the test sample in focus.
- (ii) **BLAST-partial** — The prediction score for the j -th microlabel of a particular test sample \mathbf{x}_i is calculated as the maximum sequence identity score of the test sample to all training sequences annotated with the term j [10].

Naive

Here, the relative frequency of a term in the training set is used as the prediction probability for the term in all protein sequences in the test set [10, 44].

DeepGOCNN

This utilizes a 1D convolutional neural network (CNN) to learn important features from a protein sequence. The input to the network is a one-hot representation of the amino acids in the protein’s primary sequence. It applies a linear projection layer on a series of 1D convolution operation using various filter sizes to capture relevant sub-sequences that are closely related to the function of the protein [20, 21].

DeepGOMLP

This method uses a multi-layer perceptron (MLP) network to annotate proteins. It consists of two perceptron blocks each consisting of a linear function, onto which a ReLU activation, batch normalization and dropout operations are applied in successive order. The output of the first block is connected to the output of the second block to maintain the flow of gradients during training. A sigmoid activation is finally applied to the representations learned by the feed-forward layers to produce a classification output. DeepGOMLP uses the full binary and highly sparse vector (>14k dimensions) of InterPro fingerprints as input to the network [20].

NetGO3.0

NetGO3.0 [40] is an upgraded version of state-of-the-art NetGO/NetGO2.0 [43] and GOLabeler [42] models developed in previous CAFA competitions. It consists of seven component methods: Naive, BLAST-KNN, Net-KNN, LR-3mer, LR-InterPro, LR-Text and LR-ESM. Naive assigns a term to a protein based on the empirical probability of the term in the training set. BLAST-KNN annotates a protein with a functional

term based on its top-K BLAST hits. Net-KNN assigns functional terms to a protein based on the protein’s top-K interacting proteins from its PPI network data. Different logistic regression (LR) classifiers are trained for each label in the multi-label target vector using the frequency of amino acid trigrams in the protein’s amino acid sequence (LR-3mer), InterPro fingerprints (LR-InterPro) and text curated from research and protein databases (LR-text) respectively. LR-ESM trains a LR for each functional term using sequence embeddings generated by ESM1-b [30] protein language model. The predicted scores for each microlabel in each component method are ranked and the top-k ranked terms over all component models are chosen as the final prediction.

2.7 Evaluation metrics

We used standard CAFA evaluation metrics in our experiments [10, 15, 28, 44]. Mathematical definitions for all evaluation metrics are summarised in Table 3. We assessed model performance using maximum F_1 -score (F_{max}). From the precision-recall curve at varying decision thresholds, the F_{max} is calculated as the harmonic mean of the precision and recall point that gives the highest F_1 -score. This score reflects the pronounced class-imbalance in the dataset. We also report the ability of models to correctly predict the positive terms in the test set using the area under the precision-recall (PR) curve (AUPRC) metric. The ability of models to discriminate between the positive and negative classes at varying prediction thresholds, is also assessed using the area under the receiver operating characteristics (AUROC) curve [13].

Additionally, we compared model performance based on weighted F_{max} , (WF_{max}) and minimum semantic distance (S_{min}). These metrics weight predicted terms by their information content, taking into account the hierarchical nature of the ontology. Large importance is placed on highly specific terms while little importance is given to less specific labels [10, 26, 27]. In the computation of the conditional information content for each term, we followed the true-path annotation rule and calculated each term’s empirical distribution as the relative frequency of the term in the dataset, given that its parent terms are also annotated.

3 Results

Here, we present the outcomes from the experimental validation of our model. We contrast our model’s performance with commonly adopted baselines in CAFA competitions — first, with BLAST, which leverages sequence similarity [1, 25, 32], and second, with a frequency-based approach, termed Naive [10]. Additionally, we compare our model’s predictive accuracy with state-of-the-art methods in protein function prediction — DeepGOCNN, DeepGOMLP [20, 21] and NetGO3.0 [40].

We present the experimental results on Dataset-1 and Dataset-2 in Sections 3.1 and 3.2 respectively. Considering that the time of release of NetGO3.0 overlaps with the period for the curation of Dataset-1, we do not include comparison to NetGO3.0 in the results for Dataset-1. This is due to the fact that it is only available as a webserver, hence we are unable to guarantee that the sequences in Dataset-1 set do not overlap with those used in training the NetGO3.0 model. On Dataset-2, however, we show comparison to the NetGO3.0 model.

Table 3 Mathematical definitions of evaluation metrics: Below, τ is the prediction threshold, Y_i is the groundtruth multilabel and \hat{Y}_i is the predicted multilabel at threshold τ , i.e. $\hat{Y}_i(\tau) = \mathbb{1}(\hat{Y}_i \geq \tau)$. $m(\tau)$ denotes the number of proteins in the test set for which one of the predicted scores is at least τ , n_y is the size of the test set and n_Y denotes the dimension of the multi-label target vector.

Metric	Definition
Indicator function	$\mathbb{1}_{\mathcal{X}}(x) = \begin{cases} 1 & \text{if } x \in \mathcal{X} \\ 0 & \text{if } x \notin \mathcal{X} \end{cases}$
Information content	$I(v) = -\log_2 P(v \mid \mathcal{P}a(v)),$ $\mathcal{P}a(v)$ refers to the parent(s) of term v in the ontology.
Maximum (F_{max})	$pr(\tau) = \frac{1}{m(\tau)} \sum_{i=1}^{m(\tau)} \frac{\sum_v \mathbb{1}(v \in \{\hat{Y}_i(\tau) \cap Y_i\})}{\sum_v \mathbb{1}(v \in \{\hat{Y}_i(\tau)\})}, \quad \text{precision}$ $rc(\tau) = \frac{1}{n_y} \sum_{i=1}^{n_y(\tau)} \frac{\sum_v \mathbb{1}(v \in \{\hat{Y}_i(\tau) \cap Y_i\})}{\sum_v \mathbb{1}(v \in \{Y_i\})}, \quad \text{recall}$ $F_{max} = \max_{\tau} \left\{ 2 \times \frac{pr(\tau) \times rc(\tau)}{pr(\tau) + rc(\tau)} \right\}$
Weighted (WF_{max})	$wpr(\tau) = \frac{1}{m(\tau)} \sum_{i=1}^{m(\tau)} \frac{\sum_v I(v) \cdot \mathbb{1}(v \in \{\hat{Y}_i(\tau) \cap Y_i\})}{\sum_v I(v) \cdot \mathbb{1}(v \in \{\hat{Y}_i(\tau)\})}, \quad \text{weighted precision}$ $wrc(\tau) = \frac{1}{n_y(\tau)} \sum_{i=1}^{n_y(\tau)} \frac{\sum_v I(v) \cdot \mathbb{1}(v \in \{\hat{Y}_i(\tau) \cap Y_i\})}{\sum_v I(v) \cdot \mathbb{1}(v \in \{Y_i\})}, \quad \text{weighted recall}$ $WF_{max} = \max_{\tau} \left\{ 2 \times \frac{wpr(\tau) \times wrc(\tau)}{wpr(\tau) + wrc(\tau)} \right\}$
Minimum semantic distance (S_{min})	$ru(\tau) = \frac{1}{n_y} \sum_{i=1}^{n_y} \sum_{v \in \{Y_i\} \setminus \{\hat{Y}_i(\tau)\}} I(v), \quad \text{remaining uncertainty}$ $mi(\tau) = \frac{1}{n_y} \sum_{i=1}^{n_y} \sum_{v \in \{\hat{Y}_i\} \setminus \{Y_i(\tau)\}} I(v), \quad \text{missing information}$ $S_{min}(\tau) = \min_{\tau} \left\{ \sqrt{(ru^2(\tau) + mi^2(\tau))} \right\}$
Area Under Precision Recall Curve	$AUPRC = \sum_{i=1}^{n_{\tau}} pr(\tau_i) \cdot (rc(\tau_i) - rc(\tau_{i-1})),$ n_{τ} is the number of thresholds
Area Under Receiver Operating Characteris- tics Curve	$AUROC = \sum_{i=1}^{n_{\tau}} TPR(\tau_i) \cdot (FPR(\tau_i) - FPR(\tau_{i-1})),$ TPR and FPR denote True and False positive rates respectively

3.1 Experimental results on Dataset-1

Cross-validation setup for Dataset-1

In order to reduce the risk of exaggerating generalization performance [39], we performed a homology separation between the train and test sets. We clustered the sequences in Dataset-1 at a 30% sequence identity cut-off using mmseqs2 [33]. Proteins within a cluster have at least 30% sequence similarity and 60% coverage with the cluster

Table 4 Performance evaluation based on area under precision recall curve, (AUPRC, ↑ higher the better) and maximum F_1 -score, (F_{max} , ↑ higher the better). Metrics are reported as mean ± Standard Deviations (SDs) in the form mean(SD) over 10 Cross Validation (CV) folds in all 3 ontologies. AUPRC is not reported for BLAST-full since it outputs only binary-valued predictions. A perfect prediction has $F_{max} = 1$ and AUPRC = 1. Best performing models are indicated in bold font.

Model	F_{max} (↑)			AUPRC (↑)		
	MFO	CCO	BPO	MFO	CCO	BPO
Naive	0.401(0.007)	0.619(0.007)	0.347(0.004)	0.177(0.005)	0.429(0.007)	0.231(0.006)
BLAST-partial	0.443(0.013)	0.415(0.007)	0.282(0.008)	0.284(0.009)	0.258(0.007)	0.142(0.005)
BLAST-full	0.545(0.013)	0.578(0.025)	0.357(0.012)	-	-	-
DeepGOCNN	0.501(0.011)	0.655(0.005)	0.376(0.005)	0.284(0.007)	0.207(0.011)	0.268(0.003)
DeepGOMLP	0.673(0.008)	0.657(0.006)	0.454(0.005)	0.498(0.014)	0.475(0.013)	0.397(0.007)
GO-LTR	0.682 (0.007)	0.722 (0.006)	0.486 (0.006)	0.716 (0.009)	0.787 (0.006)	0.481 (0.007)

Table 5 Performance evaluation based on weighted maximum F_1 -score, (WF_{max} , ↑ higher the better) and minimum semantic distance (S_{min} , ↓ lower the better). Metrics are reported in the form mean(SD) over 10 CV folds in all 3 ontologies. Prediction scores for terms are weighted by their conditional information content in the dataset. Thus, higher weights are given to more informative, deep and specific terms. Best performing models are highlighted in bold font.

Model	WF_{max} (↑)			S_{min} (↓)		
	MFO	CCO	BPO	MFO	CCO	BPO
Naive	0.198(0.016)	0.320(0.036)	0.258(0.009)	11.9(0.310)	10.9(0.321)	40.0(1.76)
BLAST-partial	0.383(0.017)	0.292(0.010)	0.230(0.006)	71.4(11.3)	84.3(8.78)	408.0(85.6)
BLAST-full	0.443(0.013)	0.380(0.022)	0.263(0.010)	11.5(1.09)	12.3(1.64)	46.1(1.64)
DeepGOCNN	0.344(0.015)	0.434(0.006)	0.274(0.005)	11.1(0.204)	10.1(0.210)	40.6(1.73)
DeepGOMLP	0.586(0.010)	0.464(0.006)	0.367(0.006)	8.17(0.311)	9.74(0.218)	37.7(1.15)
GO-LTR	0.591 (0.008)	0.573 (0.009)	0.392 (0.004)	7.92 (0.263)	8.16 (0.214)	34.5 (1.37)

representative (i.e centroid). This enforces a between-cluster similarity of < 30%. Clusters are iteratively refined to improve the within-cluster homology as measured by the sequence similarity. We then randomly selected 90% of the clusters for training and 10% for testing in a 10-fold cross-validation setting. All proteins in a cluster are wholly included in the train or test set. We trained models separately for each ontology. Models are optimized using 10-fold cross-validation. 10% of the sequences in the training set are used as a validation set. After the parameter optimization process for each fold, we then retrained the models on the full training set and evaluated the performance on the held-out test set.

Performance comparison of GO-LTR and competing methods

As shown in Table 4, we evaluated the predictive accuracy of our model using maximum F_1 -score (F_{max}), one of the metrics used in CAFA evaluations. Consistent with previous studies, it is evident that the machine learning (ML) methods (DeepGOCNN, DeepGOMLP and GO-LTR) outperform the common baselines used by the function

prediction community [17, 28, 44]. GO-LTR recorded the best performance compared to all competing methods in all 3 categories. In MFO, GO-LTR slightly outperformed DeepGOMLP, the second best method, with a marginal 1.3% difference in F_{max} . In predicting the cellular locations in the CCO category, GO-LTR recorded a significant performance improvement (0.718) over second-placed DeepGOMLP (0.657) model. In BPO, however, all models recorded F_{max} below 0.5, the worst performance compared to MFO and CCO function categories. This can be attributed to the dense and highly unbalanced nature of the BPO subontology. Also, the BPO graph mainly comprises broad and highly-unspecified terms. Additionally, homology-based BLAST methods were better at transferring functional terms than frequency-based Naive in MFO, but not in CCO nor BPO. Due to the highly-skewed nature of the GO dataset, we also assessed our model's performance using the area under precision recall curve (AUPRC). From Table 4, it is evident that GO-LTR's ability to predict the positive classes in the test set far exceeds that of all other competing methods. Although only modest performance differences are seen between GO-LTR and the ML baselines under the F_{max} metric, GO-LTR's strong classification performance is seen more clearly under the AUPRC metric. This indicates that GO-LTR is highly robust with respect to the choice of the prediction threshold between positive and negative classes. The Precision-Recall and ROC curves for all models are shown in Supplementary Figures A7-A8.

Information-theoretic assessment of model performance

In Table 5, we compared the predictive accuracy of models using information theoretic measures, weighted F_{max} and minimum semantic distance (S_{min}). As expected, it is seen that weighting predicted terms by their conditional information content resulted in a reduction in model performance from higher values in F_{max} in Table 4 to moderately lower values in WF_{max} in Table 5. For instance, GO-LTR's performance dropped from an F_{max} of 0.682 to a WF_{max} of 0.593. Even after the inclusion of each term's information content, GO-LTR still showed an advantage over all other methods in predicting highly specific terms in the ontology. Considering the S_{min} metric, BLAST-partial showed the least promise in identifying labels of high informative value, manifesting even worse performances in the relatively easy cases of MFO and CCO.

Contribution of different features to GO-LTR's performance

We analysed the contribution of different features to GO-LTR's predictive performance. We note that the combination of features in a multi-view learning paradigm could have complementary, redundant or contradictory effects. The results of this analyses are summarised in Table 6. The best performing GO-LTR model in MFO utilized a combination of InterPro fingerprints and sequence embeddings (UniProt). Although the third-order GO-LTR model exploiting all 3 features had the same predictive accuracy as its second-order (2-view) counterpart using InterPro and UniProt, we chose the latter model owing to its parsimonious nature. Using the network data alone (PPI) resulted in the worst performance in the MFO branch of the ontology. In the CCO category, where the goal is predicting the cellular location of proteins, it is seen that PPI had a better predictive accuracy compared to InterPro. This improved performance compared to that in MFO corroborates the assertion that proteins working together tend to be

Table 6 Ablation experiment: Effect of feature combinations on GO-LTR performance as measured by F_{max} , ↑ higher the better, in all 3 ontologies. Metric is reported as mean(SD) over 10 CV folds. Best performing feature combinations are highlighted in bold font.

Views	Feature combinations	$F_{max}(\uparrow)$		
		MFO	CCO	BPO
1-view	InterPro	0.610 _(0.010)	0.657 _(0.006)	0.409 _(0.005)
	PPI	0.491 _(0.009)	0.661 _(0.007)	0.389 _(0.006)
	UniProt	0.651 _(0.010)	0.708 _(0.005)	0.463 _(0.005)
2-view	InterPro + PPI	0.644 _(0.009)	0.682 _(0.006)	0.445 _(0.006)
	InterPro + UniProt	0.682 _(0.008)	0.710 _(0.006)	0.481 _(0.006)
	PPI + UniProt	0.670 _(0.009)	0.722 _(0.006)	0.484 _(0.006)
3-view	InterPro + PPI + UniProt	0.682 _(0.007)	0.718 _(0.006)	0.486 _(0.006)

Table 7 Ablation experiment: Performance comparison of machine learning models using all 3 features as input. Evaluation metrics are reported as mean(SD) over 10 CV folds in all 3 ontologies. Best performing models are indicated in bold font.

Model	$F_{max}(\uparrow)$			AUPRC (↑)		
	MFO	CCO	BPO	MFO	CCO	BPO
CNN-3-view	0.561 _(0.010)	0.671 _(0.006)	0.415 _(0.007)	0.314 _(0.010)	0.194 _(0.007)	0.281 _(0.006)
MLP-3-view	0.689 _(0.007)	0.726 _(0.006)	0.500 _(0.006)	0.615 _(0.010)	0.587 _(0.008)	0.487 _(0.009)
LTR-3-view	0.681 _(0.007)	0.718 _(0.006)	0.486 _(0.006)	0.710 _(0.010)	0.779 _(0.006)	0.481 _(0.007)

situated in close proximity to one another. Combination of all 3 features, however, did not yield any substantial improvement over the best performing 2-view model in the CCO category. Notably, we see that the best performing model in BPO used a combination of all three features. Indeed, all features were important in predicting terms in the BPO graph owing to its inherent complexity and dense nature. In consonance with previous works [8, 17, 28], we see that the sequence embeddings (UniProt) had the most prominent predictive signal compared to the other 2 features in all 3 ontology categories. The results in Supplementary Table A4 and Supplementary Figures A2, A3 and A4 further highlights the contributions of each feature to the predictive accuracies of the 2-view and 3-view GO-LTR models.

Performance comparison using all 3 features in ML-based models.

Next, we investigated the predictive accuracy of competing ML models by using all 3 features as input to the models. The results are presented in Table 7. In MLP-3-view, the concatenation of all 3 features was used as the new input to the original DeepGOMLP model. In respect of CNN-3-view, we concatenated the features learned by the top layers of the CNN architecture from the 1D protein sequence with the InterPro fingerprints and the PPI data. This new representation was then passed as input to the subsequent layers of the DeepGOCNN model. It is seen that the exploitation of different features enhanced the predictive accuracies of both DeepGO

models. Specifically, we see substantial improvement of $\approx 20\%$ and $\approx 5\%$ in the AUPRC of DeepGOMLP and DeepGOCNN respectively. This implies that the underlying associations between the features learned by the models resulted in improved precision and recall compared to their 1-view equivalents reported in Table 4. As shown in Table 7, MLP-3-view had the best performance in all 3 ontologies, closely followed by its GO-LTR counterpart. Although GO-LTR learns the explicit polynomial interaction between features, it achieves a similar performance to MLP-3-view which leveraged a concatenation of all 3 feature sets as input. The accompanying plots for the PR and ROC curves are shown in Supplementary Figures A9-A10.

Performance evaluation on subset of terms: depth categorizations

Further, we studied the annotation accuracy of models on different subset of labels considering their depths in the ontology. In Figure 2, we compared the performance of all models on different subset of terms differentiated by their depths in the ontology. Terms located on depths 8-11 were chosen as deep level terms, those on depths 3-7 were selected as middle level terms and labels above the third level were considered shallow level nodes. In MFO, GO-LTR outperformed all models across all 3 depth categories in the ontology (Fig 2a). In CCO, however, the Naive model had the best performance on deep level terms, closely followed by GO-LTR (Fig 2b). GO-LTR recorded the highest accuracy on the middle and shallow zones of the CCO ontology. Similarly, in BPO, Naive showed a competitive performance, even outperforming all models in predicting the shallow level nodes (Fig 2c). GO-LTR recorded the best accuracy in predicting nodes located in the deep and middle zones of the BPO category.

Performance evaluation on subset of terms: frequency categorizations

Furthermore, we evaluated the performance of all models in predicting subsets of terms grouped by their annotation frequency in the training set. Here, terms with <100 sequence examples were categorized as low frequency labels, those with frequencies in the range [100, 500) were labelled as medium frequency labels and terms having >500 examples were chosen as high frequency labels. In figure 3a, GO-LTR showed competitive performance across all 3 frequency groupings with small variations in the performance over the 10 cross validation folds. As shown in figure 3b, we see that frequency-based Naive model outperformed all models for subset of highly frequent terms in the CCO category. This means that the term frequency alone contains a high signal for predicting such terms. Hence, an appropriate combination of the predictions of machine learning and frequency-based methods could lead to performance improvement. GO-LTR exhibited the best performance on all frequency classes in the BPO ontology (Figure 3c).

Performance stability analyses based on optimal prediction threshold

Due to the extreme class-imbalance in the dataset, an adjustment to the decision threshold is necessary to reflect this bias and obtain optimal performance. As such, we assessed the stability of model predictions using the optimal prediction threshold (τ_{opt}), the threshold yielding the maximum F_1 -score. This analysis gives an overview of a model's robustness and sensitivity to small perturbations in the underlying data.

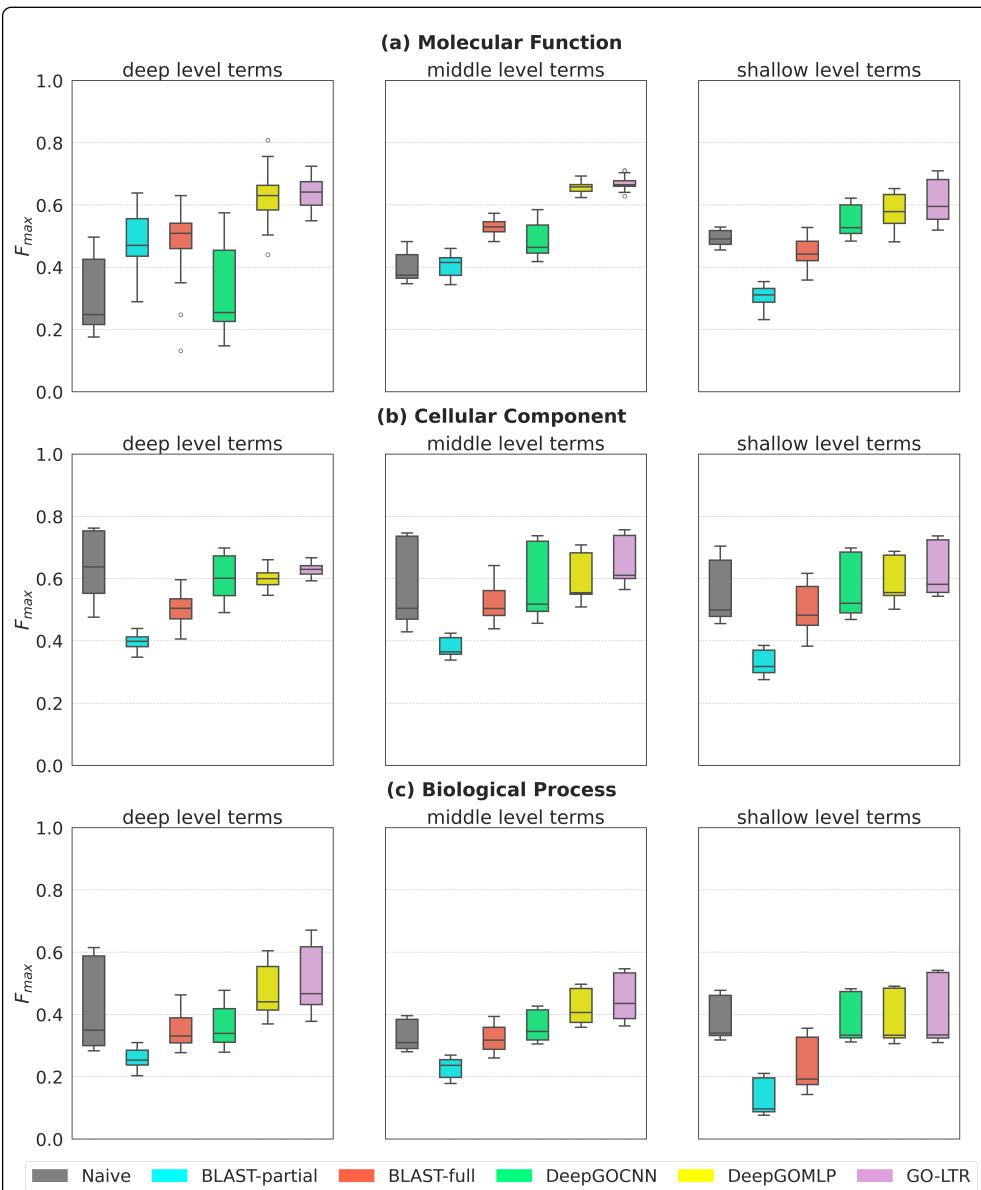
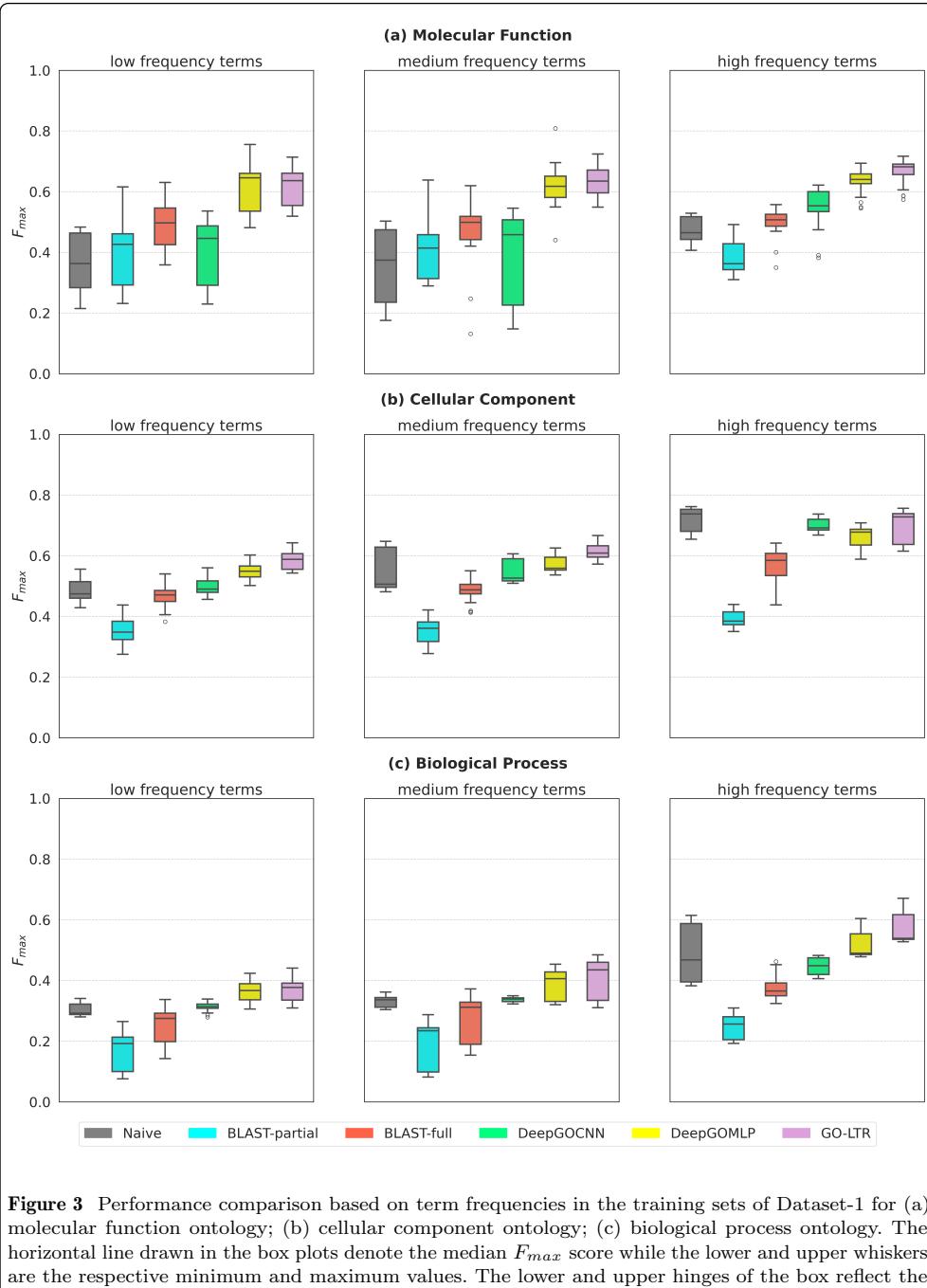
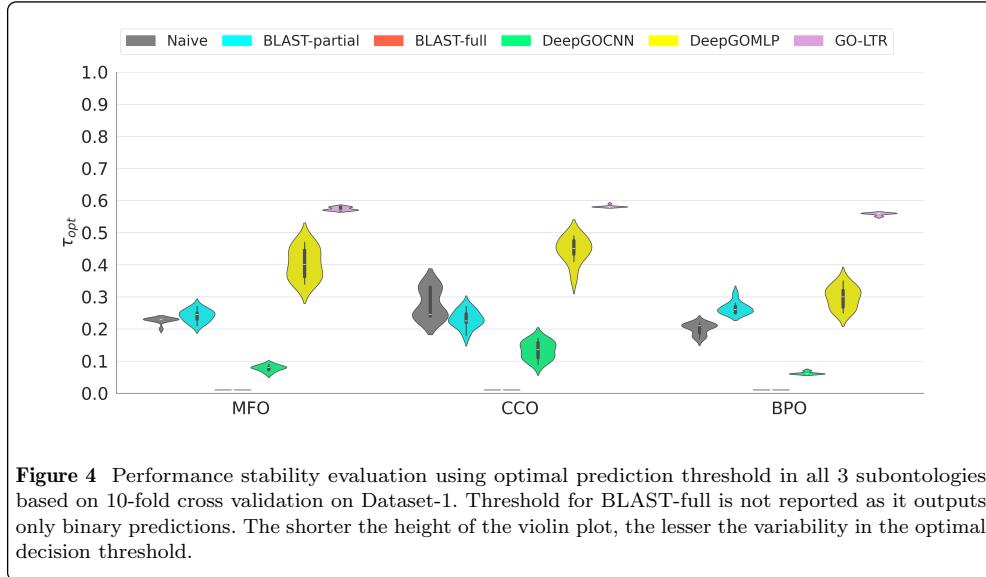


Figure 2 Performance evaluation based on depth of terms in the training sets of Dataset-1 for (a) molecular function ontology; (b) cellular component ontology; (c) biological process ontology. The horizontal line drawn in the box plots denote the median F_{max} score while the lower and upper whiskers are the respective minimum and maximum values. The lower and upper hinges of the box reflect the 25th and 75th percentiles respectively. The points outside the minimum and maximum are the outliers. The greater the median the better the performance.



As shown in Figure 4, GO-LTR exhibited a highly stable performance across the 10 cross validation (CV) folds in all 3 ontologies. DeepGO methods on the other hand recorded relatively large fluctuations in the optimal prediction threshold, making them prone to distributional shifts. We hypothesize that the huge sizes of DeepGO models as depicted by the number of trainable model parameters (Supplementary Table A5) may account for their relatively high variances compared to small-sized GO-LTR.



3.2 Experimental results on Dataset-2

Here, models are retrained on sequences in Dataset-1 and predictions are made for proteins in Dataset-2. Hence, Dataset-2 is used as an unseen dataset on which generalization performance is assessed. Models are trained separately for each ontology. Additionally, we submitted sequences in Dataset-2 to the NetGO3.0 webserver and recorded the results.

Performance evaluation on Dataset-2

We compared the generalization performance of all models on an independent test set, Dataset-2. From the precision-recall curves in Figure 5, GO-LTR achieves competitive annotation accuracy in the MFO category, performing on par with the DeepGOMLP model. In CCO, however, GO-LTR's performance surpassed all other models. In BPO, NetGO3.0, which leverages multiple features ranging from research text, term frequency, sequence embeddings, PPI, InterPro and sequence-similarity-based annotation transfer, came in first place, outperforming both DeepGOMLP and GO-LTR. The strong performance exhibited by multi-modal NetGO3.0 and the results for 3-view GO-LTR

in Supplementary Table A4 highlights the crucial importance of multi-view methods in improving annotation accuracy on the BPO category of the gene ontology. From the results of information theoretic performance evaluation presented in Supplementary Tables A1, A2 and A3, we see that GO-LTR exhibits a strong potential in predicting highly specific and rarely annotated terms in the various ontologies. In the ROC space, where a model's ability to discriminate between positive and negative classes is assessed, GO-LTR, again, shows a highly competitive generalization performance (Supplementary Figure A1).

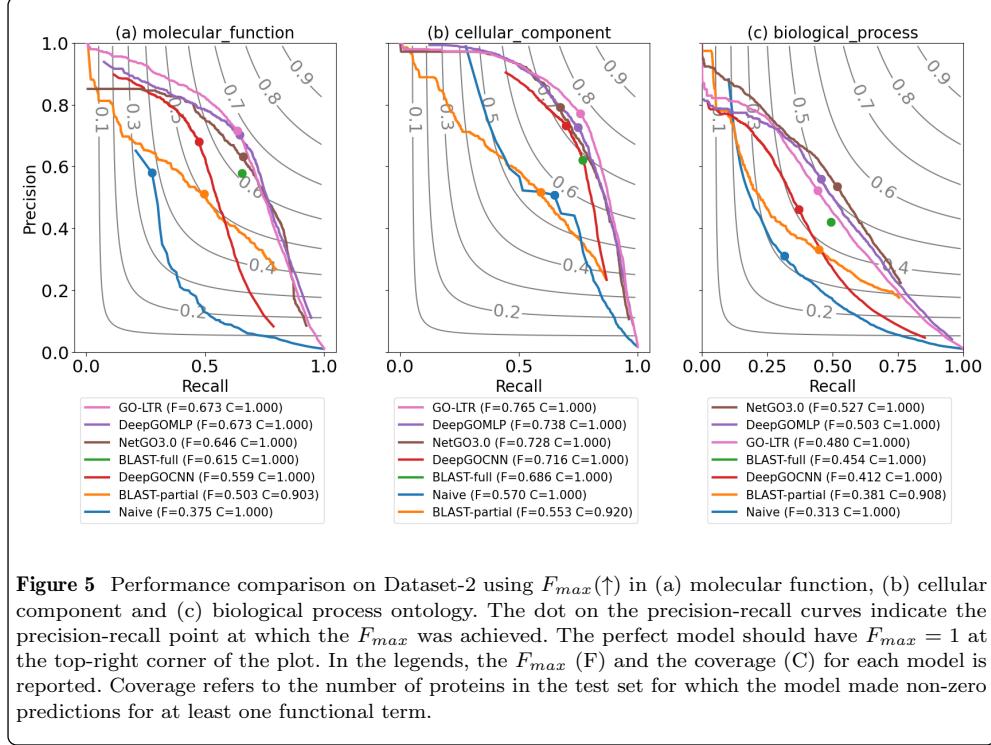
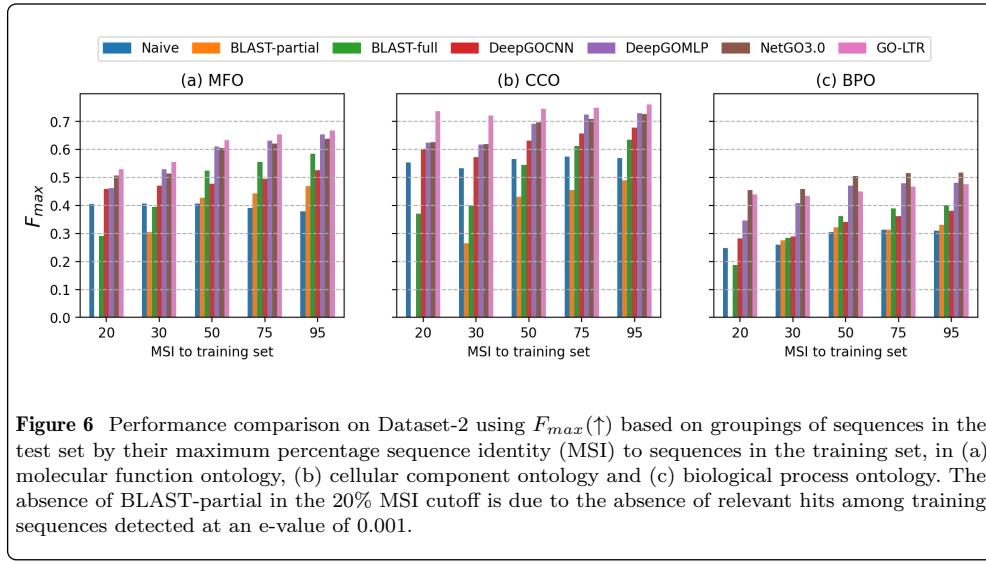


Figure 5 Performance comparison on Dataset-2 using $F_{max}(\uparrow)$ in (a) molecular function, (b) cellular component and (c) biological process ontology. The dot on the precision-recall curves indicate the precision-recall point at which the F_{max} was achieved. The perfect model should have $F_{max} = 1$ at the top-right corner of the plot. In the legends, the F_{max} (F) and the coverage (C) for each model is reported. Coverage refers to the number of proteins in the test set for which the model made non-zero predictions for at least one functional term.

Performance comparison based on maximum sequence identity

We investigated the practical utility of GO-LTR versus all competing methods in annotating novel sequences at varying degrees of homology to the sequences present in the training set. We partitioned the sequences in the test set into 5 groups ($\leq 20\%$, $\leq 30\%$, $\leq 50\%$, $\leq 75\%$ and $\leq 95\%$) based on their maximum percentage sequence identity (MSI) to those in the training set. As illustrated in Figure 6a, GO-LTR achieved the best generalization performance across all sequence similarity thresholds in the MFO category, including the highest performance on sequences with very low homology ($\leq 20\%$ and $\leq 30\%$) to known sequences in the training set. Similar results asserting the good performance of GO-LTR are shown in Figure 6b for the CCO function

category. In BPO, however, NetGO3.0 recorded the highest performance, followed closely by GO-LTR for very low sequence similarity cutoffs (Figure 6c). DeepGOMLP edged past GO-LTR slightly for sequences with relatively high homology ($> 50\%$ similarity) to the training set. As anticipated, similarity-based BLAST recorded worse performances than ML-based methods for very low sequence identity thresholds, in all ontologies, improving substantially only with an increase in maximum sequence similarity. The results of information theoretic assessments at varying identity thresholds (Supplementary Figures A5, A6) further highlight the promising potential of GO-LTR in annotation novel sequences with highly specific and rarely observed functional terms in the ontology.



4 Discussion

In this work, we studied how different protein features can be integrated for the protein function annotation task in a multi-view learning framework. Specifically, we introduced GO-LTR, a latent tensor reconstruction model that learns the multi-way interactions between multiple protein features. Extensive evaluation across several performance measures demonstrate the competitive predictive accuracy of GO-LTR in annotating proteins.

The experimental findings show that performance improvements are seen when using feature combinations in all 3 ontologies. Notably, the incorporation of all 3 features as input resulted in the best performance in the BPO category. Modest performance degradation, however, was seen when utilizing all 3 features in the molecular function and cellular component sub-ontologies. Additionally, the integration of information derived from the motifs and families in the InterPro feature and the

sequence embeddings (UniProt) resulted in a marked improvement over the use of each feature separately. These ontology-specific performance discrepancies can be attributed to the inherent intricacies of each sub-ontology.

Interestingly, even though the network data (PPI) for some proteins in our dataset were absent in the StringDB data, using the dense representation equivalent of this feature independently in CCO outperformed the case where InterPro fingerprints were used. This is likely explained by the close proximity of proteins working in concert in a living cell. Also, motifs and domain information in the InterPro fingerprints were highly predictive of functions in BPO. We posit that InterPro fingerprints are highly influential in describing larger cellular processes like those in BPO. Furthermore, the best performing linear GO-LTR model, making use of only one feature, outperformed the frequency and similarity-based baseline methods in all three sub-ontologies, thereby highlighting the practical potential of GO-LTR in annotating proteins.

Indeed, the result of further ablation studies illustrate that feature combinations are necessary in improving the generalization performance and stability of all machine learning models considered in this study. Surprisingly, while DeepGOMLP leverages residual connections between the outputs of successive layers to maintain the flow of gradients, its annotation accuracy in the 3-view experiments was not significantly better than that of GO-LTR which has no skip connections in its architecture. Additionally, the number of model parameters in the biggest GO-LTR model, the 3-view model, is several orders of magnitude less than the 1-view DeepGO counterparts.

The observations from the performance comparison based on depth and frequency of functional terms, indicate the competitive annotation accuracy of GO-LTR in predicting highly specific and rarely observed terms in the gene ontology. As the prediction of highly informative terms is desired in this task, we propose that future studies should include information content of terms with respect to the ontology, in the optimization objective.

We investigated the predictive accuracies of all competing methods on an unseen dataset. The findings show GO-LTR’s strong performance in the MFO and CCO categories of the ontology. GO-LTR, however, fell to third place in the BPO category, outperformed by NetGO3.0 and DeepGOMLP by 4.7% and 2.3% respectively. Consistent with previous works in the function prediction space, all models exhibited substantially worse performances in the BPO category compared to the results in MFO and CCO subontologies. This substantial drop in performance could be attributed to the deep and dense nature of the BPO graph, as well as the low annotation quality and high-level abstraction of its terms. Additional studies investigating this low performance phenomenon from several modeling and experimental perspectives are required.

We studied the generalizability of models to sequences at varying homologies to observed sequences in the training set. The results assert GO-LTR’s practical potential in annotating proteins that fall in the most difficult, midnight zone (< 20%) of sequence identity. This is essential because in this zone, inference via homology-based methods, which capitalize on evolutionary relatedness, provide highly unreliable and statistically uncertain results [9, 31]. Similarly, GO-LTR, like all other machine learning models outperformed homology-based BLAST and term-frequency-dependent Naive model, in the twilight (20 – 35%) and safe (> 40%) zones of sequence similarity. These results

contextualize the importance of data-dependent models and multiple informative features in reliably predicting functional terms.

5 Conclusion

In this study, we introduced GO-LTR, an automatic protein function prediction method that leverages the underlying relationships between diverse protein features in a multi-view learning framework. It relies on an efficient tensor-based estimation of model parameters. Extensive experimental validation demonstrate its high prospects in annotating proteins under several challenging conditions: generalizing to low sequence homology, rarely observed functional terms and highly specialized terms in the gene ontology.

Declarations

Supplementary information. The supporting information accompanying the paper can be found in the Appendix

Funding. We acknowledge Jane and Aatos Erkko Foundation funding under project no. 220048 (Virtual laboratory for Biodesign, JAES-BIODESIGN), as well as Research Council of Finland (grants 339421 and 345802) and the support from the Center for Young Synbio Scientists (CYSS).

Acknowledgements. Computing resources were provided by the Aalto Science-IT project.

Competing interest. None declared.

Authors' contributions. R.E.A-S., S.S., and J.R., conceived the research idea. R.E.A-S., S.S., and J.R., developed computational methods and evaluation schemes. R.E.A-S., and S.S. performed the computational experiments. R.E.A-S., S.S., and J.R., wrote the paper.

Code availability. The source code for running the model is available at <https://github.com/aalto-ics-kepaco/GO-LTR-prediction> under MIT License.

Availability of data and materials. The dataset used in the experiments and their descriptions are available at <https://github.com/aalto-ics-kepaco/GO-LTR-prediction>

Ethics approval and consent to participate. Not applicable.

Consent for publication. Not applicable

References

- [1] Altschul SF, Gish W, Miller W, et al (1990) Basic local alignment search tool. *Journal of molecular biology* 215(3):403–410
- [2] Apweiler R, Attwood TK, Bairoch A, et al (2000) Interpro—an integrated documentation resource for protein families, domains and functional sites. *Bioinformatics* 16(12):1145–1150
- [3] Ashburner M, Ball CA, Blake JA, et al (2000) Gene ontology: tool for the unification of biology. *Nature genetics* 25(1):25–29
- [4] Ashburner M, Ball CA, Blake JA, et al (2001) Creating the gene ontology resource: design and implementation. *Genome research* 11(8):1425–1433
- [5] Bairoch A, Apweiler R, Wu CH, et al (2005) The universal protein resource (uniprot). *Nucleic acids research* 33(suppl_1):D154–D159
- [6] Bateman A, Martin MJ, Orchard S, et al (2022) Uniprot: the universal protein knowledgebase in 2023. *Nucleic acids research* 51(D1)
- [7] Blondel M, Fujino A, Ueda N, et al (2016) Higher-order factorization machines. *Advances in Neural Information Processing Systems* 29
- [8] Cao Y, Shen Y (2021) Tale: Transformer-based protein function annotation with joint sequence-label embedding. *Bioinformatics* 37(18):2825–2833
- [9] Chung SY, Subbiah S (1996) A structural explanation for the twilight zone of protein sequence homology. *Structure* 4(10):1123–1127
- [10] Clark WT, Radivojac P (2013) Information-theoretic evaluation of predicted ontological annotations. *Bioinformatics* 29(13):i53–i61
- [11] Consortium GO (2023) Download the ontology. Gene Ontology resource: <https://purl.obolibrary.org/obo/go/go-basic.obo>
- [12] Database BF (2024) BFD Downloads. <https://bfd.mmseqs.com/>
- [13] Davis J, Goadrich M (2006) The relationship between precision-recall and roc curves. In: Proceedings of the 23rd international conference on Machine learning, pp 233–240
- [14] Elnaggar A, Heinzinger M, Dallago C, et al (2021) Prottrans: Toward understanding the language of life through self-supervised learning. *IEEE transactions on pattern analysis and machine intelligence* 44(10):7112–7127
- [15] Friedberg I, Radivojac P (2017) Community-wide evaluation of computational function prediction. *The Gene Ontology Handbook* pp 133–146

- [16] Gligorijević V, Renfrew PD, Kosciolek T, et al (2021) Structure-based protein function prediction using graph convolutional networks. *Nature communications* 12(1):3168
- [17] Jiang Y, Oron TR, Clark WT, et al (2016) An expanded evaluation of protein function prediction methods shows an improvement in accuracy. *Genome biology* 17(1):1–19
- [18] Johnson A, Lewis J, ALBERTS B (2002) Molecular biology of the cell. Garland Science New York
- [19] Kaltofen E, Trager BM (1990) Computing with polynomials given by black boxes for their evaluations: Greatest common divisors, factorization, separation of numerators and denominators. *Journal of Symbolic Computation* 9(3):301–320
- [20] Kulmanov M, Hoehndorf R (2022) Deepgozero: improving protein function prediction from sequence and zero-shot learning based on ontology axioms. *Bioinformatics* 38(Supplement_1):i238–i245
- [21] Kulmanov M, Khan MA, Hoehndorf R (2018) Deepgo: predicting protein functions from sequence and interactions using a deep ontology-aware classifier. *Bioinformatics* 34(4):660–668
- [22] Lewin B (2007) Cells. Jones & Bartlett Learning
- [23] Mitchell A, Chang HY, Daugherty L, et al (2015) The interpro protein families database: the classification resource after 15 years. *Nucleic acids research* 43(D1):D213–D221
- [24] Mitchell AL, Attwood TK, Babbitt PC, et al (2019) Interpro in 2019: improving coverage, classification and access to protein sequence annotations. *Nucleic acids research* 47(D1):D351–D360
- [25] Needleman SB, Wunsch CD (1970) A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of molecular biology* 48(3):443–453
- [26] Paolis CD (2023) Information Accretion. GitHub: <https://github.com/claradepaolis/InformationAccretion/tree/main>
- [27] Piovesan D, Davzago, Joshi P (2023) CAFA-evaluator. GitHub: <https://github.com/BioComputingUP/CAFA-evaluator/tree/kaggle>
- [28] Radivojac P, Clark WT, Oron TR, et al (2013) A large-scale evaluation of computational protein function prediction. *Nature methods* 10(3):221–227
- [29] Rendle S (2010) Factorization machines. In: 2010 IEEE International conference on data mining, IEEE, pp 995–1000

- [30] Rives A, Meier J, Sercu T, et al (2021) Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences* 118(15):e2016239118
- [31] Rost B (1999) Twilight zone of protein sequence alignments. *Protein engineering* 12(2):85–94
- [32] Smith TF, Waterman MS (1980) New stratigraphic correlation techniques. *The Journal of Geology* 88(4):451–457
- [33] Steinegger M, Söding J (2017) Mmseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nature biotechnology* 35(11):1026–1028
- [34] Suzek BE, Wang Y, Huang H, et al (2015) Uniref clusters: a comprehensive and scalable alternative for improving sequence similarity searches. *Bioinformatics* 31(6):926–932
- [35] Szedmak S, Cichonska A, Julkunen H, et al (2020) A solution for large scale nonlinear regression with high rank and degree at constant memory complexity via latent tensor reconstruction. *arXiv preprint arXiv:200501538*
- [36] Szklarczyk D, Morris JH, Cook H, et al (2016) The string database in 2017: quality-controlled protein–protein association networks, made broadly accessible. *Nucleic acids research* p gkw937
- [37] Vaswani A, Shazeer N, Parmar N, et al (2017) Attention is all you need. *Advances in neural information processing systems* 30
- [38] Von Mering C, Jensen LJ, Kuhn M, et al (2007) String 7—recent developments in the integration and prediction of protein interactions. *Nucleic acids research* 35(suppl.1):D358–D362
- [39] Walsh I, Pollastri G, Tosatto SC (2016) Correct machine learning on protein sequences: a peer-reviewing perspective. *Briefings in bioinformatics* 17(5):831–840
- [40] Wang S, You R, Liu Y, et al (2023) Netgo 3.0: protein language model improves large-scale functional annotations. *Genomics, Proteomics & Bioinformatics* 21(2):349–358
- [41] Wang T, Szedmak S, Wang H, et al (2021) Modeling drug combination effects via latent tensor reconstruction. *Bioinformatics* 37(Supplement_1):i93–i101
- [42] You R, Zhang Z, Xiong Y, et al (2018) Golabeler: improving sequence-based large-scale protein function prediction by learning to rank. *Bioinformatics* 34(14):2465–2473

- [43] You R, Yao S, Xiong Y, et al (2019) Netgo: improving large-scale protein function prediction with massive network information. *Nucleic acids research* 47(W1):W379–W387
- [44] Zhou N, Jiang Y, Bergquist TR, et al (2019) 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

Appendix A Protein Function Prediction Through Multi-view Multi-label Latent Tensor Reconstruction.

(by Robert Ebo Armah-Sekum, Sandor Szedmak, Juho Rousu),
Department of Computer Science, Aalto University, Finland

A.1 Performance comparison on Dataset-2 (i.e. independent test set)

Dataset-2 is described in the Section 2.3 of the main text. It comprises sequences in the UniProtKB database that gained functional annotations in the period: 14.03.2023 - 24.01.2024.

Table A1 Performance evaluation on independent test set for Molecular Function Ontology. Best performing model is highlighted in bold font and second place model is underlined.

Model	$F_{max}(\uparrow)$	$WF_{max}(\uparrow)$	$S_{min}(\downarrow)$	AUPRC(\uparrow)	AUROC(\uparrow)
Naive	0.375	0.147	10.153	0.131	0.851
BLAST-partial	0.503	0.461	13.507	0.444	0.886
BLAST-full	0.615	0.524	9.305	-	0.791
DeepGOCNN	0.559	0.417	8.085	0.396	0.906
DeepGOMLP	0.673	0.597	<u>6.529</u>	0.601	0.971
NetGO3.0	0.646	0.576	6.667	<u>0.618</u>	0.950
GO-LTR	0.673	<u>0.585</u>	6.179	0.689	<u>0.968</u>

Table A2 Performance evaluation on independent test set for Cellular Component Ontology. Best performing model is highlighted in bold font and second place model is underlined.

Model	$F_{max}(\uparrow)$	$WF_{max}(\uparrow)$	$S_{min}(\downarrow)$	AUPRC(\uparrow)	AUROC(\uparrow)
Naive	0.570	0.439	7.461	0.327	0.929
BLAST-partial	0.553	0.481	14.626	0.522	0.898
BLAST-full	0.686	0.562	9.219	-	0.842
DeepGOCNN	0.716	0.584	5.927	0.301	0.939
DeepGOMLP	<u>0.738</u>	<u>0.622</u>	<u>5.816</u>	0.681	0.975
NetGO3.0	0.728	0.611	6.300	<u>0.781</u>	0.958
GO-LTR	0.765	0.663	5.142	0.817	0.975

Table A3 Performance evaluation on independent test set for Biological Process Ontology. Best performing model is highlighted in bold font and second place model is underlined.

Model	$F_{max}(\uparrow)$	$WF_{max}(\uparrow)$	$S_{min}(\downarrow)$	AUPRC(\uparrow)	AUROC(\uparrow)
Naive	0.313	0.218	27.437	0.181	0.872
BLAST-partial	0.381	0.340	87.963	0.317	0.844
BLAST-full	0.454	0.387	32.526	-	0.730
DeepGOCNN	0.412	0.309	26.846	0.331	0.898
DeepGOMLP	<u>0.503</u>	<u>0.423</u>	22.965	<u>0.463</u>	0.945
NetGO3.0	0.527	0.454	21.916	0.484	0.823
GO-LTR	0.480	0.386	22.696	0.455	<u>0.942</u>

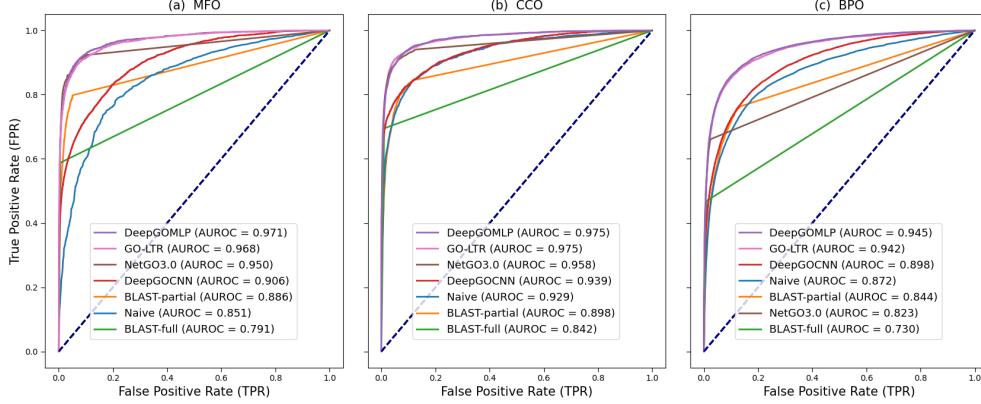


Figure A1 Performance comparison on Dataset-2 using receiver operating characteristics (ROC) curves in (a) molecular function, (b) cellular component and (c) biological process ontology. The area under the ROC (AUROC) curve is given in the legends.

A.1.1 Performance evaluation of feature combinations in GO-LTR on independent test set

Here, the model is retrained on the full training set (sequences annotated on or before 13.03.2023) and evaluated on the independent test set (sequences annotated between 14.03.2023 and 24.01.2024)

Table A4 Ablation experiment on independent test set: Effect of feature combinations to GO-LTR performance as measured by F_{max} , ↑ higher the better, in all 3 ontologies. Best performing feature combinations are highlighted in bold font.

Views	Feature combinations	$F_{max}(\uparrow)$		
		MFO	CCO	BPO
1-view	InterPro	0.598	0.685	0.397
	PPI	0.418	0.596	0.334
	UniProt	0.615	0.748	0.453
2-view	InterPro + PPI	0.641	0.693	0.427
	InterPro + UniProt	0.673	0.764	0.476
	PPI + UniProt	0.654	0.765	0.473
3-view	InterPro + PPI + UniProt	0.663	0.757	0.480

A.1.2 Feature contributions in GO-LTR extracted from the trained model

The effect of each feature/view d is given by $\mathbf{D}_{\lambda_U}^{(d)}$ in Table 1f. The weights corresponding to the contribution of each feature in the prediction of the output is then extracted from the retrained model described in Section A.1.1. The top-100 components are then sorted in descending order of magnitude ($\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{100}$) and their magnitude plotted on a logarithmic scale. The components are L2-normalized.

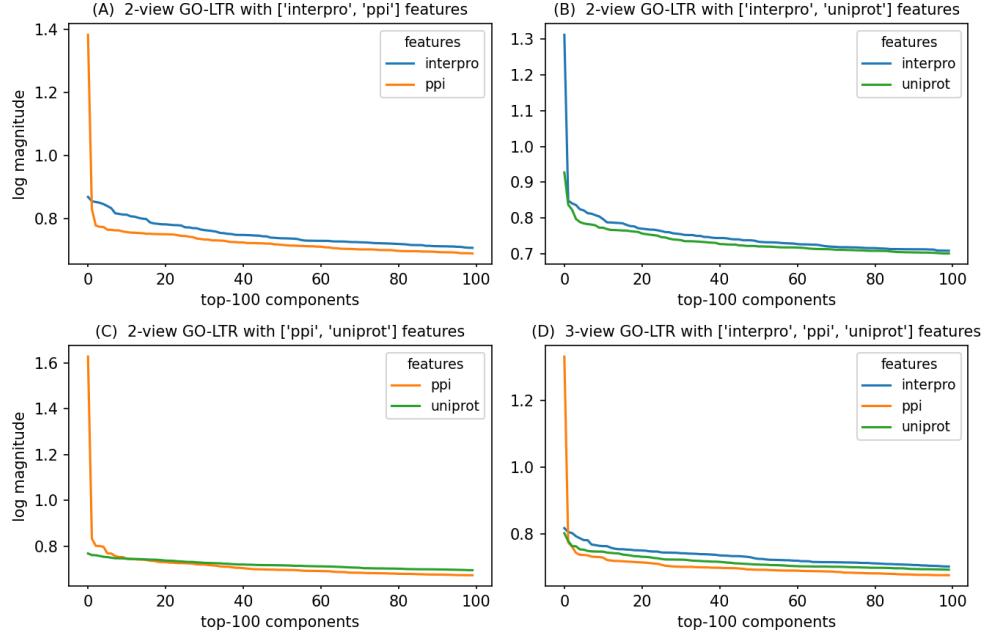


Figure A2 Feature contributions in 2-view and 3-view GO-LTR models of the Molecular Function ontology based on Dataset-2: The magnitudes of the top-100 components are plotted on a logarithmic scale.

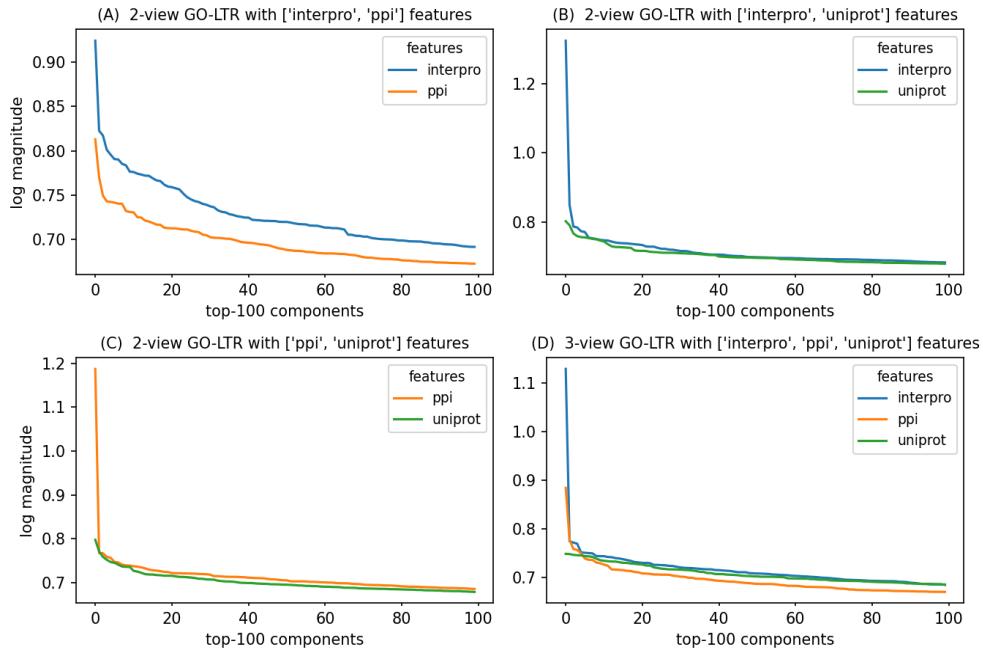


Figure A3 Feature contributions in 2-view and 3-view GO-LTR models of the Cellular Component ontology based on Dataset-2: The magnitudes of the top-100 components are plotted on a logarithmic scale.

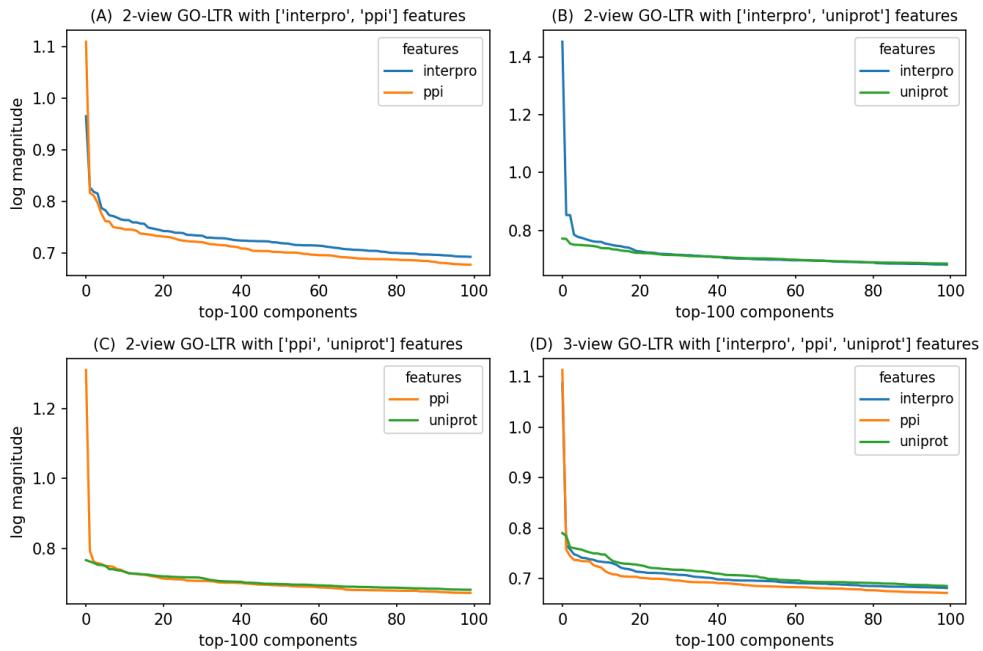


Figure A4 Feature contributions in 2-view and 3-view GO-LTR models of the Biological Process ontology based on Dataset-2: The magnitudes of the top-100 components are plotted on a logarithmic scale.

A.1.3 Performance evaluation on Dataset-2

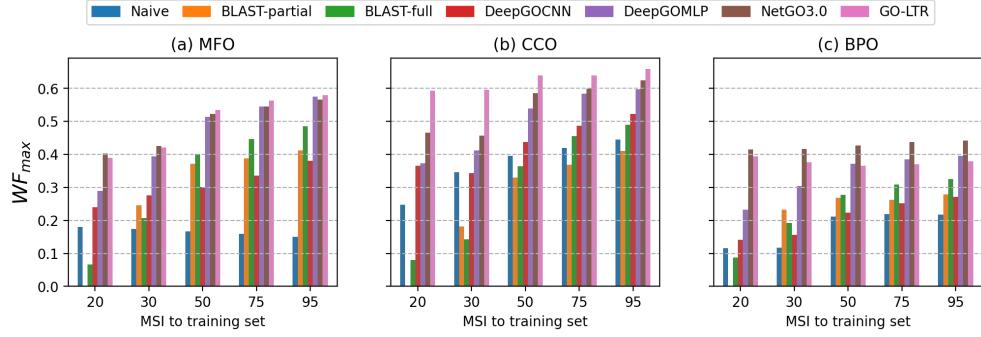


Figure A5 Performance comparison on Dataset-2 using $WF_{max}(\uparrow)$ based on groupings of sequences in the test set by their maximum percentage sequence identity (MSI) to sequences in the training set, for all ontologies. The absence of BLAST-partial in the 20% MSI cutoff is due to the absence of relevant hits among training sequences detected at an e-value of 0.001.

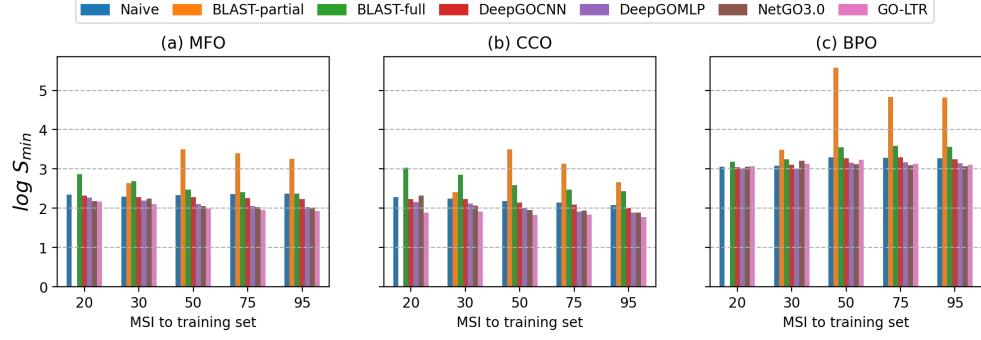


Figure A6 Performance comparison on Dataset-2 using $S_{min}(\downarrow)$ based on groupings of sequences in the test set by their maximum percentage sequence identity (MSI) to sequences in the training set, for all ontologies. The absence of BLAST-partial in the 20% MSI cutoff is due to the absence of relevant hits among training sequences detected at an e-value of 0.001.

A.2 Performance comparison on Dataset-1

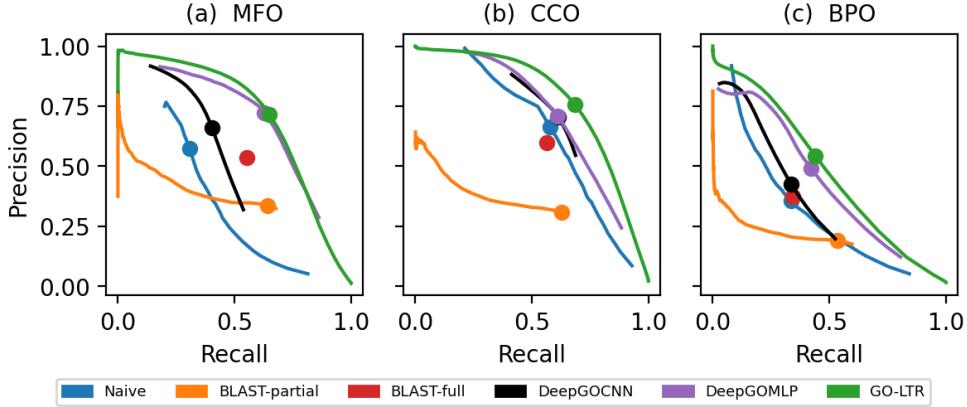


Figure A7 Mean precision-recall curves over the 10 outer cross validation splits of Dataset-1 in all ontologies. The dot on each curve indicate the precision-recall point at which the mean F_{max} was achieved.

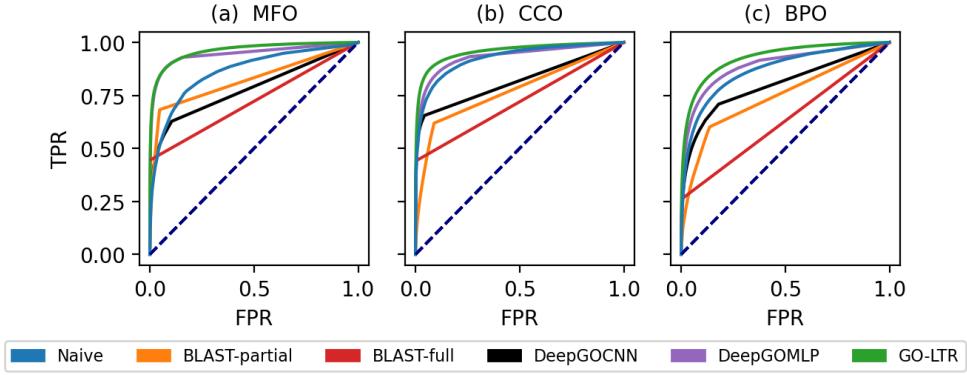


Figure A8 Mean receiver operating characteristics (ROC) curves over the 10 outer cross validation splits of Dataset-1 in all ontologies. In the figure, the mean true positive rate (TPR) versus the mean false positive rate (FPR) is plotted.

A.3 Performance comparison of 3-view Machine Learning (ML) models on Dataset-1 10-fold CV

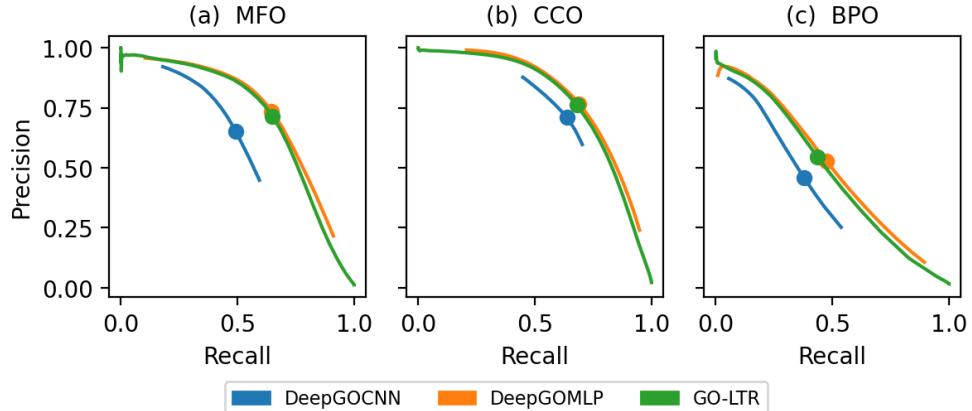


Figure A9 Mean precision-recall curves over the 10 outer cross validation splits of Dataset-1 for 3-view ML models in all ontologies. The dot on each curve indicate the precision-recall point at which the average F_{max} was achieved.

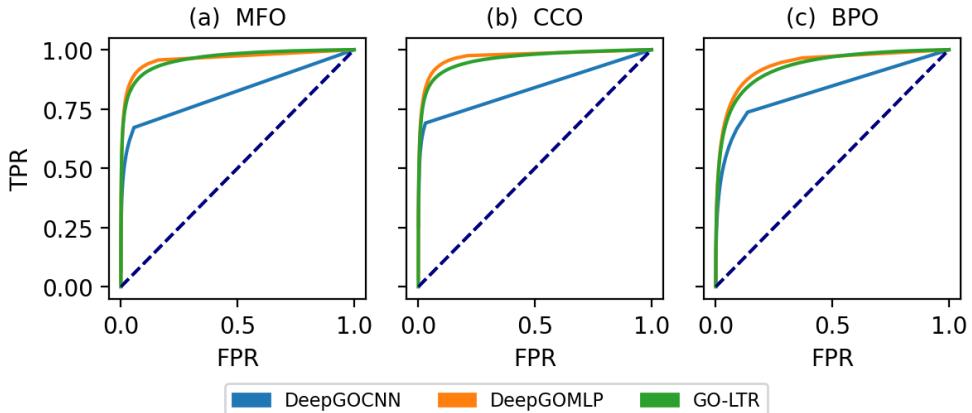


Figure A10 Mean receiver operating characteristics (ROC) curves over the 10 outer cross validation splits of Dataset-1 for 3-view ML models in all ontologies. In the figure, the mean true positive rate (TPR) versus the mean false positive rate (FPR) is plotted.

A.4 Comparison of model parameters in ML models

Table A5 Number of trainable parameters in DeepGOCNN, DeepGOMLP and 3-view GO-LTR model. Models with the least number of parameters are indicated in bold font.

Model	MFO	CCO	BPO
DeepGOCNN	20,891,400	20,726,375	23,221,225
DeepGOMLP	17,192,712	17,975,911	20,756,457
3-view GO-LTR	3,614,400	3,614,400	3,614,400

A.5 Architecture of ProtT5 protein language model

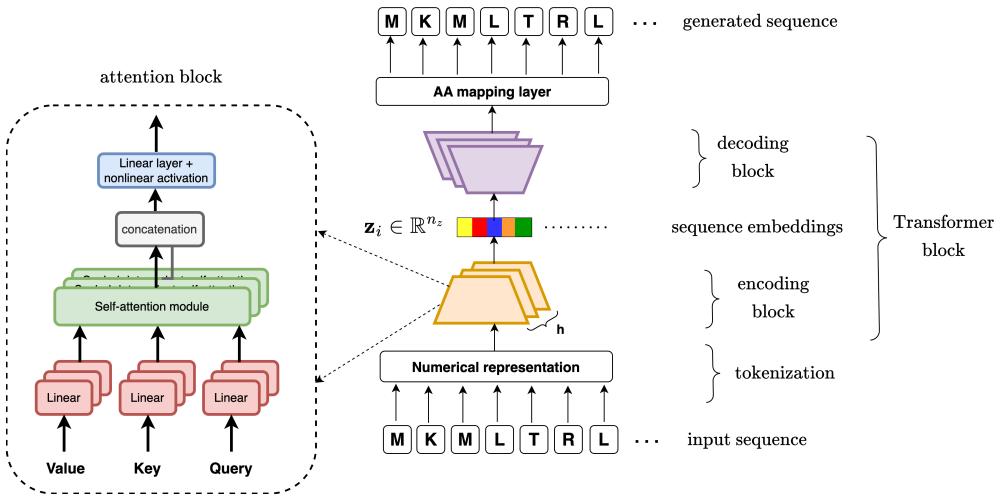


Figure A11 Encoding block: In the architecture, the input AA sequence is converted to a numerical representation (tokenization) for subsequent layers. The encoding block consists of several encoding units stacked on top of one another, with each learning complex relationships from the input sequence. Each encoding unit is made up of a self-attention layer and a feed-forward network (FFN). A Layer Normalization operation is computed atop the output of each attention and FFN layer. The attention layer helps to capture dependencies between each AA and the remaining residues. There are multiple attention units (i.e. attention heads) running in parallel to learn diverse relationships in the input sequence. The FFN module learns a non-linear projection of the learned relationships into a low-dimensional space (embeddings). **Decoding block:** In a similar fashion, the decoding block comprises several decoding units stacked on top of one another, with the output of one feeding into the input for the next. The decoder generates the full-length protein sequence in an auto-regressive manner: using masked self-attention on the right-shifted output sequence, cross attention on the representations from the encoding block to generate a residue at a time by paying attention to only the residues observed up to the current time point. Just like that observed in the encoding block, there are multiple attention heads running in parallel to learn richer and diverse representations encoded by the previous layer during the sequence generation process. Finally, the numerical representations of the generated residues are then mapped back into the AA character space. We would like to note that the embeddings utilized in our studies are extracted from the encoder's output.