

## METHOD

# QAUST: Protein Function Prediction Using Structure Similarity, Protein Interaction, and Functional Motifs



Fatima Zohra Smali<sup>1,#</sup>, Shuye Tian<sup>2,#</sup>, Ambrish Roy<sup>3</sup>, Meshari Alazmi<sup>1,4</sup>,  
 Stefan T. Arold<sup>5</sup>, Srayanta Mukherjee<sup>3</sup>, P. Scott Hefty<sup>6</sup>, Wei Chen<sup>2,\*</sup>, Xin Gao<sup>1,\*</sup>

<sup>1</sup>Computational Bioscience Research Center (CBRC), Computer, Electrical and Mathematical Sciences and Engineering (CEMSE) Division, King Abdullah University of Science and Technology (KAUST), Thuwal 23955, Saudi Arabia

<sup>2</sup>Department of Biology, Southern University of Science and Technology of China (SUSTC), Shenzhen 518055, China

<sup>3</sup>Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI 48109, USA

<sup>4</sup>College of Computer Science and Engineering, University of Ha'il, Ha'il 55476, Saudi Arabia

<sup>5</sup>Biological and Environmental Sciences and Engineering (BESE) Division, King Abdullah University of Science and Technology (KAUST), Thuwal 23955, Saudi Arabia

<sup>6</sup>Department of Molecular Bioscience, University of Kansas, Lawrence, KS 66047, USA

Received 11 July 2018; revised 3 April 2019; accepted 17 May 2019  
 Available online 23 February 2021

Handled by Ziding Zhang

**Abstract** The number of available protein sequences in public databases is increasing exponentially. However, a significant percentage of these sequences lack functional annotation, which is essential for the understanding of how biological systems operate. Here, we propose a novel method, Quantitative Annotation of Unknown Structure (QAUST), to infer protein functions, specifically Gene Ontology (GO) terms and Enzyme Commission (EC) numbers. QAUST uses three sources of information: structure information encoded by global and local structure similarity search, biological network information inferred by protein–protein interaction data, and sequence information extracted from functionally discriminative sequence motifs. These three pieces of information are combined by consensus averaging to make the final prediction. Our approach has been tested on 500 protein targets from the Critical Assessment of Functional Annotation (CAFA) benchmark set. The results show that our method provides accurate functional annotation and outperforms other prediction methods based on sequence similarity search or threading. We further demonstrate that a previously unknown function of human tripartite motif-containing 22 (TRIM22) protein predicted by QAUST can be experimentally validated.

**KEYWORDS** Protein function prediction; GO term; EC number; Protein structure similarity; Functionally discriminative motif

## Introduction

As of today, over 150 million protein sequences are available in the UniProtKB/TrEMBL database [1]. However,

this increase in the number of known protein sequences does not reflect a parallel increase in our biological knowledge, as less than 1% of these sequences have a manually annotated function [2]. On the other hand, the functional annotation of these sequences is not only an essential step for the understanding of physiological processes and biological systems in living entities, but also one of the

\*Corresponding authors.

E-mail: [xin.gao@kaust.edu.sa](mailto:xin.gao@kaust.edu.sa) (Gao X), [chenw@sustech.edu.cn](mailto:chenw@sustech.edu.cn) (Chen W).

#Equal contribution.

Peer review under responsibility of Beijing Institute of Genomics, Chinese Academy of Sciences / China National Center for Bioinformation and Genetics Society of China. <https://doi.org/10.1016/j.gpb.2021.02.001>

1672-0229 © 2021 The Authors. Published by Elsevier B.V. and Science Press on behalf of Beijing Institute of Genomics, Chinese Academy of Sciences / China National Center for Bioinformation and Genetics Society of China.

This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

highly challenging tasks in biology, which is why there is an increasing need to provide reliable, automated protein function annotation.

Significant efforts have been made to identify evolutionarily related proteins and automatically transfer functional annotations between homologous protein pairs [3–6]. To make such sequence similarity-based functional transfer possible, powerful sequence-alignment methodologies have been developed. In particular, algorithms like BLAST/PSI-BLAST [3] and hidden Markov model (HMM)-based techniques [4–6] have been frequently used to transfer functional annotations between homologous proteins. The underlying assumption of these sequence-based methods is that evolutionarily related proteins may inherit the function of a shared common ancestor. However, there are numerous cases in which proteins with high sequence similarity have distinct functions [7,8]. To partially address the problem, several methods have been developed to predict function using annotated conserved sequence motifs that are responsible for the functional aspect of the protein. These methods typically construct the sequence motifs from multiple sequence alignment of proteins belonging to the same protein family with known function [9–11]. They, however, have two major limitations. First, high-quality sequence alignment is typically required for motif construction, which is not trivial to obtain especially when the sequence homology is low. Second, the accuracy is limited by the quality of functional annotation of motifs. To overcome these limitations, we propose in this work to use a protein-specific “functionally discriminative motif” constructed from sequence fragments excised from the template sequence.

From another perspective, the 3D structure of a protein sequence is believed to be more involved in its biological function [12,13], since structures are more conserved than the sequences. The 3D structure of a protein can therefore provide additional information for function transfer, especially when the sequence similarity between related proteins is too low for sequence homolog detection [14,15]. However, the relationship between the protein function and its structure is not straightforward, as in some cases, similar structures perform the same function while in many cases similar folds perform different functions [16,17]. Therefore, many prediction methods have been relying on local structure similarity search methods rather than global similarity search to identify functionally homologous proteins [18–20]. Most of these approaches scan the query protein against a library of known conserved spatial motifs or known active sites (*e.g.*, binding sites) with known function [21]. Local similarity search methods have been proven to be quite accurate in detecting functional similarity between proteins of different folds, but they also have a high probability of producing false positive matches [22]. One

possible solution is to combine global and local structure alignments to overcome the promiscuity of global structure comparison and low specificity of local structure matching [23,24], which we implement in this project.

A number of function prediction methods are based on the information extracted from protein–protein interaction (PPI) networks [25,26]. The assumption in this case is that proteins that physically interact with each other frequently appear at the same sub-cellular location and are part of the same biological process [27]. However, it is not always the case that proteins which interact with each other share the same molecular function [*e.g.*, Programmed cell death protein 1 (PD1) and Programmed death-ligand 1 (PD-L1)], which is why PPI information is not always sufficient to predict very specific functions [28].

Finally, recently there is an emergence of methods which combine multiple sources of information (PPI, domains, sequence alignments, *etc.*) using advanced machine learning algorithms to perform function prediction. These methods have shown improved prediction performance over methods that use only one type of information [29–37].

In this work, we propose a new protein function prediction method, Quantitative Annotation of Unknown Structure (QAUST), which combines the global and local structure similarity search with PPI networks and functional sequence motif detection. Our approach follows a sequence-to-structure-to-function workflow. Starting from the protein amino acid sequence, we first generate structure prediction by the Iterative Threading ASSEMBLY Refinement (I-TASSER) method [38]. The predicted structure is then used to identify the proteins with similar functions based on a combination of global and local structure similarity search method that follows the same pipeline used in COFACTOR [24,39]. PPI information is meanwhile extracted from the STRING database [40]. And finally, we extract functionally discriminative sequence motifs as our third main prediction feature. The confidence scores obtained from these three features are combined in a consensus function to obtain our final confidence score.

Since the terminology of a “protein function” might be ambiguous, we would like to clarify that the definitions of function followed in this work are Enzyme Commission (EC) numbers [41] and Gene Ontology (GO) terms [42]. EC numbers are used to categorize enzymes into hierarchical families using a numerical classification. Specifically, the EC number (which is composed of four numbers separated by periods, *i.e.*, A.B.C.D) refers to the reaction catalyzed by a specific enzyme. On the other hand, the GO terms are a set of controlled vocabulary to formally describe proteins and RNAs based upon their functions. Three aspects of ontologies, biological process (BP), cellular component (CC), and molecular function (MF), are defined in this database. Each one of these three GO aspects is represented by a

structured directed acyclic graph (DAG), where nodes represent GO terms that describe gene product functions, while the edges represent the relationships (“is\_a” or “part\_of”) between the GO terms. In GO functional hierarchy, the more general functions are on the top of the graph, while more specific terms are usually present further down the graph.

Our prediction results are compared to the following programs: 1) COFACTOR [39], a global and local structure similarity-based method; 2) LOMETS [43], a meta-threading algorithm; 3) HHsearch [5], an HMM-based method that is widely used to detect protein homologs; 4) BLAST [3], which transfers annotations based on sequence similarity; 5) naïve baseline, which predicts GO terms based on their annotation frequency; and 6) two highly-ranked methods from the Critical Assessment of Functional Annotation (CAFA) assessment [44], GoFDR [32] and INGA [45].

## Method

### Dataset

To evaluate QAUST for EC prediction, we use the benchmark dataset of COFACTOR [24,39] as our testing dataset. This dataset consists of 318 enzymes with unique EC numbers (first three digits) covering all 6 enzyme classes. Similarly, all sequences in our template libraries having a sequence identity > 30% with the query enzymes are excluded from the template libraries.

We evaluate QAUST for GO prediction on a dataset of 500 randomly chosen non-redundant proteins from the CAFA 2 Targets (<https://biofunctionprediction.org/cafa/>) annotated with at least one GO term. To eliminate any structure or function homologs to the query, templates having a sequence identity > 30% with the query proteins are excluded from the template libraries both in the I-TASSER threading library and our function prediction template libraries.

### EC number prediction

#### *Global and local similarity search*

The first step of our protein function prediction is the generation of the predicted 3D model of the query protein using I-TASSER [38], as outlined in “Section 1” in File S1. The predicted model of the query protein obtained from I-TASSER is then scanned against a non-redundant (pairwise sequence identity no more than 90%) structure template library of 2385 enzymes with at least the first three digits of EC number annotated by the Catalytic Site Atlas (CSA) database [46]. This library scanning detects homologous structure templates to the query proteins using

two types of structure similarity search programs: global similarity search and local similarity search.

#### Global similarity search

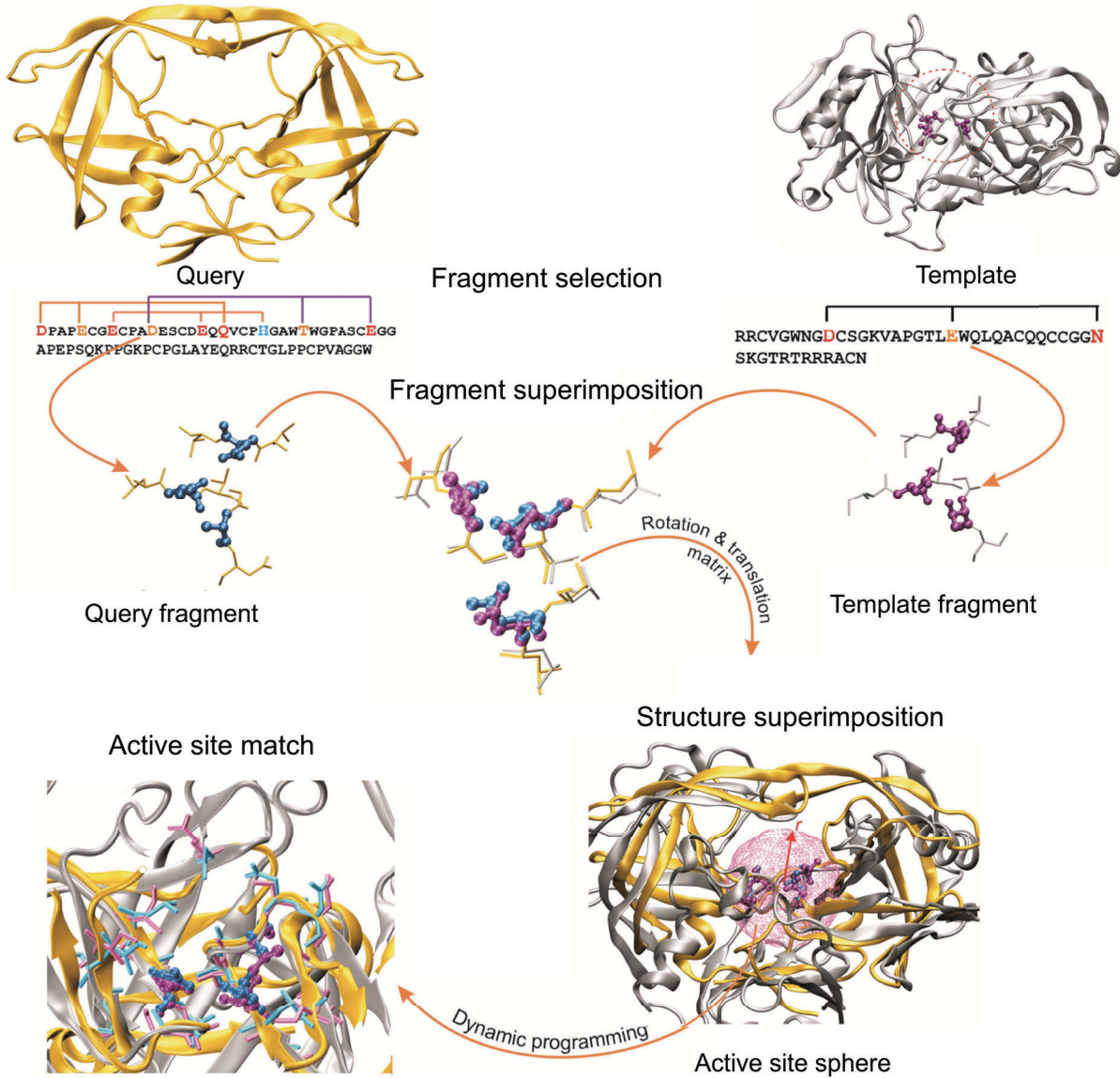
Templates with a similar global structure to the predicted structure of the query protein are detected from the template library using TM-align [47]. Another important consideration when searching for templates with similar global folds to the query protein is the quality of the structural models. Appraising the accuracy of the structure modeling in the scoring scheme helps to reduce the number of false positive predictions. In this particular case, the quality of the predicted I-TASSER model generated in the previous step is evaluated using *Cscore* [38].

#### Local similarity search

The local structural search approach consists of three steps (Figure 1). The first step is the structural match of the specific catalytic/active residue pairs. For a given pair of query and template proteins, we first scan the known catalytic/active residues of the template through the query sequence. The query’s residues whose amino acid types are the same as the amino acid types of the template’s catalytic/active residues are marked as potential active sites in the query. The structures of all combined sets of marked residues in the query are extracted from the predicted model and used as candidate active sites. The structures of the candidate sites are superimposed on the known catalytic/active residues in the template. To make the structure superimposition more reliable, for each residue  $i$ , the coordinates of  $C_{\alpha}$  atoms and side-chain centers of mass of the two neighboring residues, *i.e.*, the  $(i-1)$ -th and  $(i+1)$ -th residues, are also included in the superimposition.

The second step is to identify the key local environment residues around the active sites in the query and the template. For this purpose, we superimpose the complete structure of the query and template proteins based on the rotation matrix obtained from the superimposition of the candidate catalytic/active residue structures obtained in the previous step. A sphere of radius  $r$  is then defined around the geometric center of the template’s local 3D fragments, where  $r$  is the maximum distance of the template residues in the local 3D fragment from the geometric center. The sphere represents a local environment or a probable active site region, under which the chemical and structural similarity of the query and template are compared. Because a sphere comprising of a very small number of catalytic/active residues can easily generate false positive hits, when the template’s active site region is small, we set the number of residues inside the sphere to be a minimum of 20 residues. This value is obtained using minimum grid search parameter optimization by evaluating different sphere sizes in the range of [10, 50] residues to select the most accurate value.

In the third step, the best alignment of the local active site



**Figure 1** A schematic diagram of the local similarity search procedure for functional site identification

The residues of the query protein (yellow) in the active site region are shown in cyan, while those of the template protein (grey) are shown in magenta.

residues in the spheres between the query and the template is identified using a scoring function similar to TM-align. Starting from the initial superposition of the query and template protein structures, we perform a Needleman-Wunsch dynamic programming to generate the best alignment for the residues in the selected sphere of the template and the query, where the alignment score matrix  $S_{ij}$  for aligning the  $i$ -th residue in the query and the  $j$ -th residue in the template is defined as:

$$S_{ij} = \left[ \frac{1}{1 + \frac{d_{ij}}{d_0}} + M_{ij} \right] \quad (1)$$

where  $d_{ij}$  is the  $C_\alpha$  distance between residues  $i$  and  $j$ ,  $d_0$  is the distance cutoff given by  $d_0 = 1.24\sqrt[3]{L-15} - 1.8$  obtained

from TM-align,  $M_{ij}$  is the substitution score between the  $i$ -th and  $j$ -th residues taken from the BLOSUM62 mutation matrix with the value normalized by the diagonal element in the mutation matrix. The gap penalty is set as  $-1$ . For a given scoring matrix  $S_{ij}$ , a new alignment is generated by dynamic programming. A new superposition and scoring matrix are then constructed based on the new alignment to obtain a newer alignment from dynamic programming. This procedure is iteratively repeated until the final alignment is converged. For each alignment, the active site match (AcM) is evaluated using an alignment score defined as:

$$\text{AcM} = \frac{1}{N_t} \sum_{i=1}^{N_{ali}} \frac{1}{1 + \left(\frac{d_{ii}}{d_0}\right)^2} + \frac{1}{N_t} \sum_{i=1}^{N_{ali}} M_{ii} \quad (2)$$



where  $N_i$  represents the number of residues in the active site sphere of the template,  $N_{ali}$  is the number of aligned residue pairs. The maximum AcM score obtained during the heuristic iterations is recorded for each candidate active site. Finally, the set of residues in the candidate active site which has the highest AcM score is selected to evaluate the similarity between the active sites of the query and the template. The weights that form the AcM score have been derived based on the predicted structures of 100 randomly chosen training proteins from the template library, which are non-homologous (sequence similarity < 30%) to the test proteins in order to maximize the sensitivity and specificity of the predictions.

### Scoring function for global and local similarity search

The final score for predicting EC numbers, used to sort the hits from the enzyme library, is a combination of the global similarity search score and the AcM score (obtained from the local similarity search) and is defined as:

$$\underline{QAUSTEC} = C_{norm} \cdot \left[ TM + \frac{Cov}{1 + RMSD_{ali}} \right] + 2 \cdot ID_{ali} \cdot Cov \quad (3)$$

$$+ \frac{AcM}{2}$$

where  $Cov$  represents the coverage of the structural alignment,  $RMSD_{ali}$  is the root-mean-square deviation (RMSD) between the model and the template structure in the structurally aligned region, and  $ID_{ali}$  is the sequence identity between query and template based on the alignment generated by TM-align. The hyperbolic-tangent-like normalization is further used to normalize the raw EC score to be between 0 and 1:

$$\underline{QAUSTEC}_{norm} = \frac{2}{1 + \exp(-\underline{QAUSTEC})} - 1 \quad (4)$$

### GO term prediction

For GO term prediction, we combine three different predictors. Each one of these predictors generates a confidence score. The three confidence scores obtained are then combined in a consensus function to generate the final prediction score. The first predictor is the global structure similarity, which uses I-TASSER to predict the 3D structure of the query, and then scans a library of templates to identify those which have a similar global structure to the predicted model. The second predictor is based on PPI information, and the third one is based on extracted functional sequence motifs.

#### Global protein structure similarity

Similar to EC prediction, I-TASSER is also used here to construct the corresponding 3D model to the query sequence. The model obtained is then scanned against a library of templates to identify those which share a similar global

structure to the query model (<https://zhanglab.ccmb.med.umich.edu/BioLiP/library.html>). For the time being, the functionally important residues for most of the proteins in the GO template library are unknown. Therefore, only the global similarity search is taken into consideration when sorting the hits from the GO library. Global similarity search for GO prediction is done in a similar way to global similarity search for EC prediction described in the previous section. The only difference is that to select the best hits for GO prediction, we rank a template using the  $Fh_{score}$  defined as:

$$Fh_{score} = C_{norm} \left( TM_{score} + \frac{1}{1 + RMSD_{ali}} \cdot Cov \right) \quad (5)$$

$$+ 3 \cdot ID_{ali} \cdot Cov$$

Since each single protein can be annotated with multiple GO terms and the global search may result in many close template structures, a query protein can have multiple GO term predictions with high  $Fh_{score}$ . Therefore, the confidence score of each GO term is calculated as follows:

$$P_{structure}(\lambda) = \frac{1}{N} \sum_{i=1}^{N_i} Fh(i) \quad (6)$$

where  $\lambda$  represents a given GO term,  $N_i$  is the number of templates annotated with the GO term  $\lambda$ , and  $N$  is the total number of templates selected for generating the consensus. When multiple close templates are available, we only consider the templates with  $Fh_{score} > 1$ . For those query proteins with less than 10 templates of  $Fh_{score} > 1$ , the top 10 templates are selected for generating the consensus prediction regardless of the  $Fh_{score}$ . Also, given the hierarchical nature of the GO DAG, we consider that when a protein is annotated with a given GO term, all its ancestor GO terms (through “is\_a” relation) are automatically implied. Therefore, once a GO term  $\lambda$  is scored, we score all its ancestor terms as well. The score of any ancestor GO term  $\mu$  of term  $\lambda$  is calculated as:

$$P_{structure}(\mu) = P_{structure}(\lambda) \left( 1 + \frac{N_\mu}{N_0} \right) \quad (7)$$

where  $N_\mu$  and  $N_0$  are the numbers of leaf nodes under node  $\mu$  and the root node, respectively. Since COFACTOR [24,39] uses a similar structure scoring function, we have highlighted the main differences between QAUST and COFACTOR in “Section 2” in File S1, and compared our method with COFACTOR in the experiments.

#### PPI network

We exploit the information provided by the STRING [40] database, which is a library of PPI networks, to extend our prediction set. The query protein sequence is mapped to its corresponding STRING entry by BLAST, with minimum sequence identity cutoff of 90%. Extracting the PPI partners

of the query, we calculate the confidence score of STRING for a GO term  $\lambda$  [ $P_{STRING}(\lambda)$ ] as the frequency of the GO term  $\lambda$  among the experimentally annotated interaction partners of the query protein:

$$P_{STRING}(\lambda) = \frac{n_\lambda}{N} \quad (8)$$

where  $n_\lambda$  is the number of interaction partners annotated with the GO term  $\lambda$ , and  $N$  is the number of partners associated with term  $\lambda$ , according to the corresponding UniProt-GOA (<http://www.ebi.ac.uk/GOA>) entry of this PPI partner. This score could take any value from 0 to 1.

### Functionally discriminative sequence motifs

In addition to the structure similarity search and PPI features discussed above, we also include sequentially extracted features to predict GO terms since a sequence is a highly valuable source of information that can especially be useful when dealing with proteins for which we cannot construct a good quality 3D structure model or which with no known PPI information.

Our functionally discriminative motif detection algorithm follows three steps: detection of sequence templates, identification of functionally discriminative motifs given a GO term, and scoring the query protein.

#### Detection of sequence templates for a query protein

The sequence homologs of a query sequence are detected by PSI-BLAST [3] from the Uniref90 database [49]. We filter all obtained homologs having a sequence identity > 30% with the query.

#### Identification of functionally discriminative motifs given a GO term

We map all the selected sequence homologs of the query to their corresponding GO annotations in the UniProt-GOA database (<http://www.ebi.ac.uk/GOA>). GO terms assigned with “Inferred from Electronic Annotation (IEA)” or “No biological Data available (ND)” evidence codes are not considered. We also filter out annotations with the evidence code “Inferred from Physical Interactions (IPI)”, since we use PPI information in our features. After filtering these annotations, we are left with the annotations based on evidence codes: Inferred from Experiment (EXP), Inferred from Direct Assay (IDA), Inferred from Mutant Phenotype (IMP), Inferred from Genetic Interaction (IGI), Inferred from Expression Pattern (IEP), Traceable Author Statement (TAS), and Inferred by Curator (IC). For each GO term  $\lambda$ , we build two sets of sequences from the set of homolog sequences detected in the previous step. These two sets are: the “annotated set”, which is the set of sequence homologs annotated with this specific GO term, and the “not-annotated set”, which is the set of sequence homologs not annotated with this given GO term. For each of these two

sets, we extract the ten most frequent motifs by extracting all unique amino acid motifs of length [4, 7] from the sequence set using sliding windows. These motifs are ranked in descending order by their occurrences. The top 10 most frequent motifs are the initial “frequent list”, while the remaining motifs are in the “waiting list”. If, within the “frequent list”, a short motif is a substring of another longer motif, the shorter motif is discarded, and the most frequent motif from the “waiting list” is transferred to “frequent list” to ensure that the latter always has 10 motifs. This process is iterated until, in the “frequent list”, any motif is not a substring of another motif. The motifs in the “frequent list” are used for matching the query in the next step.

#### Scoring the query protein

For each of the two sets (annotated and not-annotated sets), we check the number of frequent motifs extracted in the previous step that are also present in the query sequence. Then, we calculate the confidence score of the GO term given the query sequence as follows:

$$P_{MOTIF}(\lambda) = \frac{n_q(\lambda)}{N(\lambda)} \left[ 1 - \frac{n_q(\lambda^c)}{N(\lambda^c)} \right] \quad (9)$$

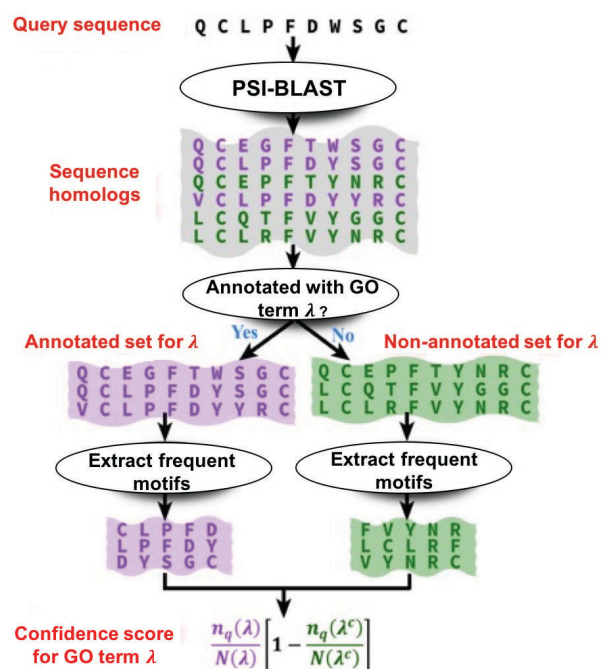
where  $\lambda$  is the given GO term.  $N(\lambda)$  and  $N(\lambda^c)$  are the numbers of frequent patterns from the “annotated set” and “not-annotated set”, respectively, both of which equal to 10.  $n(\lambda)$  and  $n(\lambda^c)$  are the corresponding numbers of matched patterns at the query sequence. This score can take any value from 0 to 1. An ideal value of this score would be equal to 1, which happens when all the sequences in the annotated set contain these frequent motifs and none of the sequences in the not-annotated set contains these same motifs. This scoring function has been designed to penalize the prediction in case the query sequence matches a high number of frequent motifs from the not-annotated set. This way, the scoring function accounts for two essential pieces of information: which set has the maximum number of frequent motifs matched in the query, and how significant is the difference between the number of matched motifs from the annotated set and that from the not-annotated set. **Figure 2** shows a flowchart detailing the three steps of extracting functional sequence motifs.

#### Consensus

To predict GO terms, the three main scores obtained from the three different predictors (the structure search, the PPI network, and the functional motifs) are combined by consensus averaging to calculate the final confidence score  $P_{consensus}(\lambda)$  for a GO term  $\lambda$ :

$$P_{consensus}(\lambda) = 1 - \prod_{m \in \{structure, STRING, MOTIF\}} (1 - P_m(\lambda)) \quad (10)$$

This equation used to calculate the consensus has been



**Figure 2** Workflow for sequence motif-based function prediction in QAUST

The query sequence is searched against the UniRef90 database [48] by PSI-BLAST to identify sequence homologs with GO term annotation. For a GO term of interest,  $\lambda$ , the identified homologs are divided into two sets: the “annotated set” (purple) which contains homologs annotated with  $\lambda$ , and the “not-annotated set” (green) which consists of homologs not associated with  $\lambda$ . From each of the two sets, frequent motifs, *i.e.*, continuous sequence fragments, are extracted. For illustration purposes, only three five-residue-long motifs from each set are drawn. The GO term  $\lambda$  is predicted with confidence score  $\frac{n_q(\lambda)}{N(\lambda)} \cdot \left[ 1 - \frac{n_q(\lambda^c)}{N(\lambda^c)} \right]$ . Here,  $N(\lambda)$  and  $N(\lambda^c)$  are the total numbers of extracted frequent motifs from “annotated set” and “not-annotated set”, correspondingly; while  $n_q(\lambda)$  and  $n_q(\lambda^c)$  are the numbers of frequent motifs from “annotated set” and “not-annotated set” that match the query sequence, respectively. In this example, only the motif “CLPFD” from “annotated set” matches the query, making the confidence score equals to  $1/3 \cdot [1 - 0/3] = 1/3$ . GO, Gene Ontology.

previously used by other methods for protein function prediction [45]. If one or more predictors are not available for a given term (*e.g.*, no interaction partners are known for the given query), only the available predictors are used to obtain the confidence score. Also, since GO uses the true-path rule (*i.e.*, if a protein is associated by a term, it is also implicitly annotated by its ancestors), for every predicted GO term, all its ancestors are considered to be predicted as well since they are more general terms.

### Cell culture, plasmid construction, and transfection

The Hek293T cells were obtained from the American Type Culture Collection (ATCC), and were cultured in DMEM (Catalog No. C11965500BT, Gibco) supplemented with 10% fetal bovine serum (FBS; Catalog No. 10270106,

Gibco) and 1% penicillin/streptomycin (P/S; Catalog No. 15070063, Gibco) with 5% CO<sub>2</sub> at 37 °C.

To establish the constructs expressing human tripartite motif-containing 22 (TRIM22) protein, we cloned the coding sequence of human TRIM22 with FLAG or GFP at its N-terminus into pcDNA3.1(+) vector, using one-step clone kit (Vazyme).

For transfection, PEI reagent was applied with the general ratio of 30 reagents as 10 μg plasmids (5 μg FLAG-TRIM22 and 5 μg GFP-TRIM22) into each 10 cm plate at 70% cell confluence. After 48 h transfection, cells were harvested for the following assays.

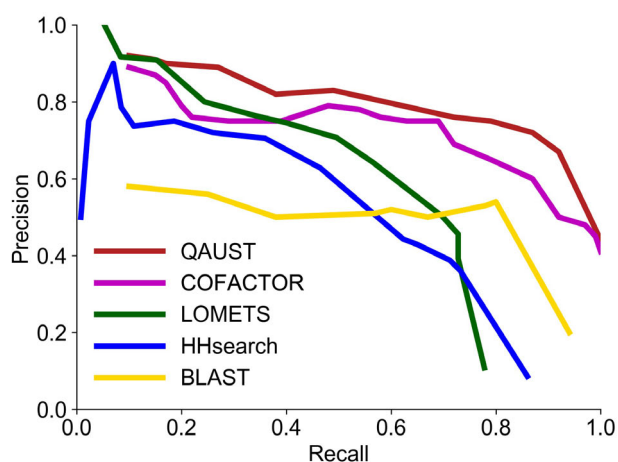
### Co-immunoprecipitation

Cells were harvested with lysis buffer (20 mM Tris-HCl pH7.5, 150 mM NaCl, 1 mM EDTA, 1% NP-40) with proteinase inhibitor. For each sample, 25 μl protein A/G beads (MCE) were incubated with 1 μg anti-FLAG (Sigma) or anti-GFP (Proteintech) antibody at 4 °C for 2 h. Mouse IgG (Cell Signaling Technology) was used as a negative control. Beads were washed three times with lysis buffer and then incubated with 500 μg cell lysate at 4 °C for another 2 h. After washed with lysis buffer three times, 20 μl 2× sample buffer were added into beads and heated at 95 °C for 10 min. Protein levels of beads or cell lysate were then detected by Western blot using individual antibodies.

## Results

### Prediction of enzyme functions (EC numbers)

We compared the EC prediction performance of our method to five other methods: HHsearch [5], LOMETS [43], BLAST [3], COFACTOR [24,39], and DEEPre webserver [49]. We compared the performance of these methods based on precision (positive predictive value) and recall (sensitivity) rates. Figure 3 shows the precision-recall graph corresponding to four baseline methods as well as QAUST. Since the DEEPre webserver does not report the confidence score with the annotation, we could not draw the precision-recall curve but compared QAUST to DEEPre based on accuracy. An EC number prediction is considered to be “true” if the first three digits of the EC number from the hit are identical to those of the query protein; otherwise the hit is considered to be “false”. As shown in Figure 3, the rate of true positive predictions using the EC-score is much higher than that of HHsearch, LOMETS, BLAST, and COFACTOR at most recall rates. QAUST has also an area under precision-recall curve (AUPRC) of 0.712 which is higher than that of COFACTOR (0.643), LOMETS (0.510), and HHsearch (0.489). Table 1 reports the accuracy of



**Figure 3** Precision-recall curves for EC prediction by QAUST, COFACTOR, LOMETS, HHsearch, and BLAST  
EC, Enzyme Commission.

**Table 1** Accuracy values of EC prediction for QAUST and five other methods

Method	Accuracy
QAUST	0.709
COFACTOR	0.698
LOMETS	0.661
HHsearch	0.607
BLAST	0.571
DEEPre	<b>0.714</b>

Note: The highest accuracy value of EC prediction is shown in bold. EC, Enzyme Commission.

QAUST compared to five other methods. The results show that DEEPre has a slightly higher performance than QAUST in terms of accuracy, which is probably due to the fact that DEEPre is a machine learning method trained on a large number of enzymes with known functions that overlap or contain close homologs to our test data.

## Prediction of GO terms

To assess the contribution of individual predictors to the GO prediction performance by QAUST, we visualized the precision-recall curve of the structure similarity search alone ( $P_{structure}$ ; corresponding to the COFACTOR method [39]), the precision-recall curve of structure similarity search combined with PPI information ( $P_{structure}$  and  $P_{STRING}$ ), and that of the final QAUST prediction ( $P_{structure}$ ,  $P_{STRING}$ , and  $P_{MOTIF}$ ). Additionally, we compared the prediction performances of our method based on different sets of features on our dataset (see the subsection Dataset in Method) to those of naïve baseline (a method predicting GO terms based on their annotation frequency), BLAST [3], LOMETS [43], HHsearch [5], INGA webserver [45] (a method combining BLAST, PPI information, and Pfam in one predictor), and GoFDR [32] (one of the top function prediction methods at

the CAFA assessment [44] which uses a machine learning model as classifier and discriminative residues as the main feature).

The performance was primarily evaluated using precision-recall curves computed at each prediction score threshold. We also used the  $F_{max}$  measure as a quantitative measure to evaluate the overall performance of the precision-recall curves. Precision, recall, and  $F_{max}$  are defined in the same way as the CAFA evaluation [50]. The  $F_{max}$  measure is computed as the maximum value of the  $F_{measure}$  which is computed at each threshold as  $\frac{2 \times precision \times recall}{precision + recall}$ .

Precision at threshold  $t$  is defined as  $\frac{|P_x(t) \cap C_x|}{|P_x(t)|}$ , while

recall is defined as  $\frac{|P_x(t) \cap C_x|}{|C_x|}$ , where  $x$  is a query protein,

$P_x(t)$  is the set of predicted terms for  $x$  at threshold  $t$ , and  $C_x$  is the set of correct terms that  $x$  is experimentally annotated with.

Similar to the CAFA evaluation [34], we also reported the minimum semantic distance ( $S_{min}$ ) as an additional evaluation metric for GO prediction.  $S_{min}$  is defined as  $min_t \left\{ \sqrt{ru(t)^2 + mi(t)^2} \right\}$ .  $ru(t)$  is the remaining uncertainty

at threshold  $t$  defined as  $\frac{1}{n_e} \sum_{i=1}^{n_e} \sum_f ic(f) | f \notin P_i(t) \wedge f \in C_i$ ,

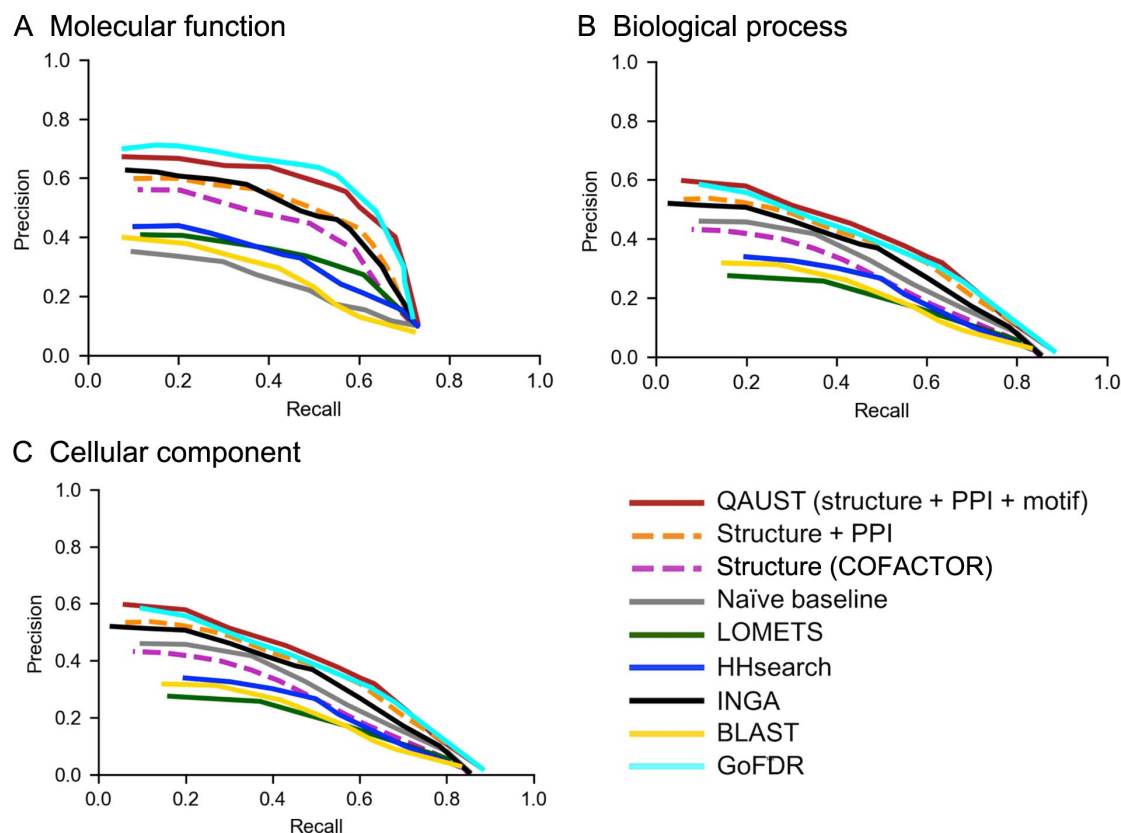
and  $mi(t)$  is the misinformation at threshold  $t$  defined as

$\frac{1}{n_e} \sum_{i=1}^{n_e} \sum_f ic(f) | f \in P_i(t) \wedge f \notin C_i$ , where  $n_e$  is the number

of proteins in our dataset,  $P_i(t)$  is the set of predicted GO terms for protein  $i$  at threshold  $t$ ,  $C_i$  is the set of terms that protein  $i$  is actually annotated with, and  $ic(f)$  is the information content of the GO term  $f$ . The very general and unspecific GO terms such as “MF”, “BP”, “CC”, “Binding”, and “Protein Binding” are excluded from the evaluation.

As shown in Figure 4, our method combining structure, PPI, and functional motif information achieves higher precision than most other methods at most recall points, in particular for BP and CC. For our dataset, structure and motif information has been used for all proteins. However, the PPI information from STRING is missing for 74 proteins. In this case, only the structure and motif information are used. The  $F_{max}$  measure values are shown in Table 2, and the  $S_{min}$  values are shown in Table 3. Surprisingly, in CC prediction, naïve baseline, which predicts GO terms based on their annotation frequency, achieves higher performance than all other methods including QAUST. In fact, in the CAFA assessment [44], naïve baseline also outperforms most of the other methods in predicting CC terms. One possible explanation for why naïve baseline has higher  $F_{max}$





**Figure 4** Precision-recall curves for GO prediction

GO prediction performances of our method based on different sets of features, and six other methods for each of the three GO branches are compared. **A.** Molecular function. **B.** Biological process. **C.** Cellular component.

**Table 2**  $F_{\max}$  values of each branch of GO for QAUST and other prediction methods

Method	MF	CC	BP
Structure + PPI + Motifs (QAUST)	0.568	0.467	<b>0.448</b>
Structure + PPI	0.507	0.453	0.436
Structure (COFACTOR)	0.467	0.402	0.367
Naïve Baseline	0.315	<b>0.492</b>	0.387
LOMETS	0.396	0.374	0.303
HHsearch	0.381	0.356	0.347
INGA	0.501	0.436	0.421
BLAST	0.347	0.373	0.321
GoFDR	<b>0.579</b>	0.449	0.431

Note: The highest  $F_{\max}$  value of each branch of GO is shown in bold. GO, Gene Ontology; MF, molecular function; CC, cellular component; BP, biological process.

**Table 3**  $S_{\min}$  values of each branch of GO for QAUST and other prediction methods

Method	MF	CC	BP
Structure + PPI + Motifs (QAUST)	7.66	5.41	<b>10.80</b>
Structure (COFACTOR)	7.51	5.81	11.72
Naïve Baseline	8.26	<b>5.12</b>	12.09
LOMETS	8.11	7.23	14.56
HHsearch	8.33	7.68	14.20
INGA	7.95	6.74	12.27
BLAST	8.42	6.34	14.01
GoFDR	<b>7.32</b>	5.60	11.65

Note: The highest  $S_{\min}$  value of each branch of GO is shown in bold.  $S_{\min}$ , minimum semantic distance.

for CC term prediction is because the most frequently used CC terms in protein annotation are usually part of a small set of very general terms such as “cytoplasm” or “intracellular part”. Since the naïve baseline is solely based on frequency, it increases the chance of predicting a true positive [44]. We also reported the  $P$  values obtained from the Mann-Whitney  $U$  test to assess the significance of the difference in performance of QAUST compared to all other methods in “Section 3” in [File S1](#) and [Table S1](#).

As a further analysis, to investigate if the performance of our method is solely due to the power of the I-TASSER structure prediction we used, we replaced the I-TASSER structure prediction component of our method by HHsearch and LOMETS structure prediction, respectively. Our results show that no matter which structure prediction method is used, our scoring function,  $P_{structure}$ , can significantly improve the performance on predicting the GO terms. Meanwhile, among the three structure prediction methods, I-TASSER with  $P_{structure}$  consistently performs the best over all three GO hierarchy branches of MF, BP, and CC, whereas LOMETS with  $P_{structure}$  has the second best performance on CC, and HHsearch with  $P_{structure}$  is the second best on predicting MF and BP terms ([File S1](#), Section 4; [Figure S1](#)). Additionally, we evaluated the performance of our method when only PPI and motif information are used without including any structure-based information. The results show that the function prediction performance drops when structure features are not used ([File S1](#), Section 5; [Figure S2](#)).

#### *How do PPI information and functional sequence motifs improve the prediction ?*

PPI information extracted from STRING is an important feature used in our prediction. In [Figure 4](#), we show how PPI information alone improves the performance achieved by the structure similarity search (orange dash lines *versus* magenta dash lines). The precision-recall curves in [Figure 4](#) show that the contribution of PPI information from STRING is very significant for CC and BP terms, especially for large recall rates. Moreover, the precision-recall curves confirm our initial hypothesis on the utility of PPI information for function annotation. As shown in [Figure 4](#), while there is some improvement in predicting MF terms, this improvement is not substantial. The reason why PPI is not particularly helpful in MF term prediction is most probably because proteins that interact with each other do not necessarily share the same specific molecular function, even when they are part of the same biological process.

In addition to the structure similarity search and the PPI features, the results show that the functional motifs extracted improve the performance of the prediction significantly. As a further analysis, we evaluated the performance of our functional motif detection method when

both predicted and experimentally annotated GO terms are taken into consideration instead of considering experimental annotations only. The results of this experiment are reported in “Section 6” in [File S1](#) and [Table S2](#). In addition to comparing the performance of our method to BLAST, LOMETS, HHsearch, and COFACTOR, we also compared it to INGA and GoFDR, two top methods from CAFA [44] in particular for MF and BP term prediction, and to naïve baseline which is one of the performance references used in CAFA.

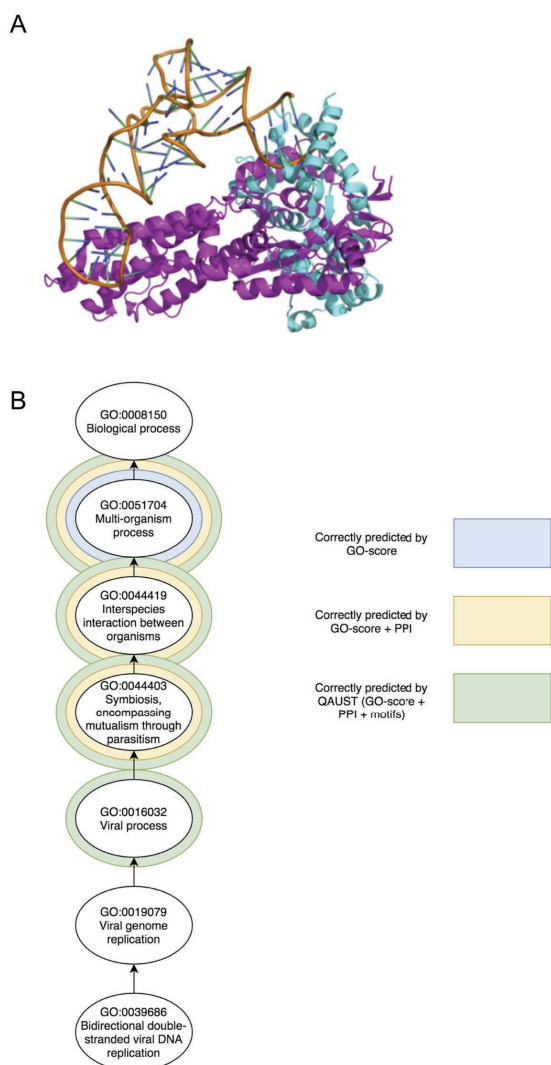
#### **Case studies**

##### *Prediction of the bacteriophage T4 gene 59 helicase assembly protein*

To better illustrate the performance of QAUST and the contribution of each component to the prediction, we used an example bacteriophage T4 gene 59 helicase assembly protein (PDB ID P13342; the cyan structure in [Figure 5A](#)), which is a DNA binding protein required mainly for DNA replication in the late stage of T4 infection [51]. [Figure 5B](#) shows the set of BP terms associated with this protein. In this particular example, both BLAST and INGA did not predict any correct term for this protein (the naïve root term is not counted here). When solely using global structure similarity ( $P_{structure}$ ), we could only predict one single correct BP term. This makes sense because all the queries in our test set are difficult targets, which do not have close homologs in the template database. For instance, the closest template for this query P13342 is the methionine-tRNA ligase (PDB ID 2CT8A; the magenta structure in [Figure 5A](#)). The sequence identity between P13342 and 2CT8A is only 6.84% and the TM-score between the two structures is only 0.24. Therefore, structure similarity or homology-based methods are not expected to predict the function of the query well. Structure information ( $P_{structure}$ ) combined with PPI predicted three correct terms out of the six experimentally annotated terms. On the other hand, QAUST predicted four correct terms out of six. In addition, the prediction of QAUST is at least one level deeper in the GO hierarchy than the other methods. Meanwhile, the predicted MF and CC terms for this protein by QAUST are at least as accurate as other methods. Two other case studies illustrating examples for predicting MF and CC terms are also shown ([File S1](#), Section 7 and Section 8; [Figures S3](#) and [S4](#)).

##### *Experimental validation of TRIM22 dimerization*

To provide an experimental assessment of the performance of QAUST, we chose the human TRIM22 protein as an example. TRIM22 is known as an interferon-inducible protein which shows antiviral activity, such as HIV, HBV, and HCV [52–54]. Recent studies have also shown that



**Figure 5** A study case for protein function prediction using QAUST **A.** The superimposition between the query (PDB ID P13342, in cyan) and the closest template in the database (PDB ID 2CT8A, in magenta) based on the structural alignment generated by TM-align. **B.** Predicted BP terms for protein with PDB ID P13342. The six BP terms (the root term, Biological Process, is a naïve term, which is not counted) shown are the experimentally annotated terms. The colored contours represent the BP terms that are predicted by the corresponding methods.

TRIM22 mediates autophagy in human macrophages [55]. However, the function of TRIM22 is still not comprehensively understood as the protein only exists in primates.

We applied QAUST to predict the function for TRIM22. Among the predicted GO terms with high consensus scores (File S1, Section 9; Table S3), some of the CC and BP terms agree well with the previously known functions of TRIM22, such as the CC term “nucleus” and the BP term “response to virus”. However, the only two predicted MF terms have quite high consensus scores, “protein binding” and “protein homodimerization activity”, suggesting that TRIM22 binds to itself to form a dimer.

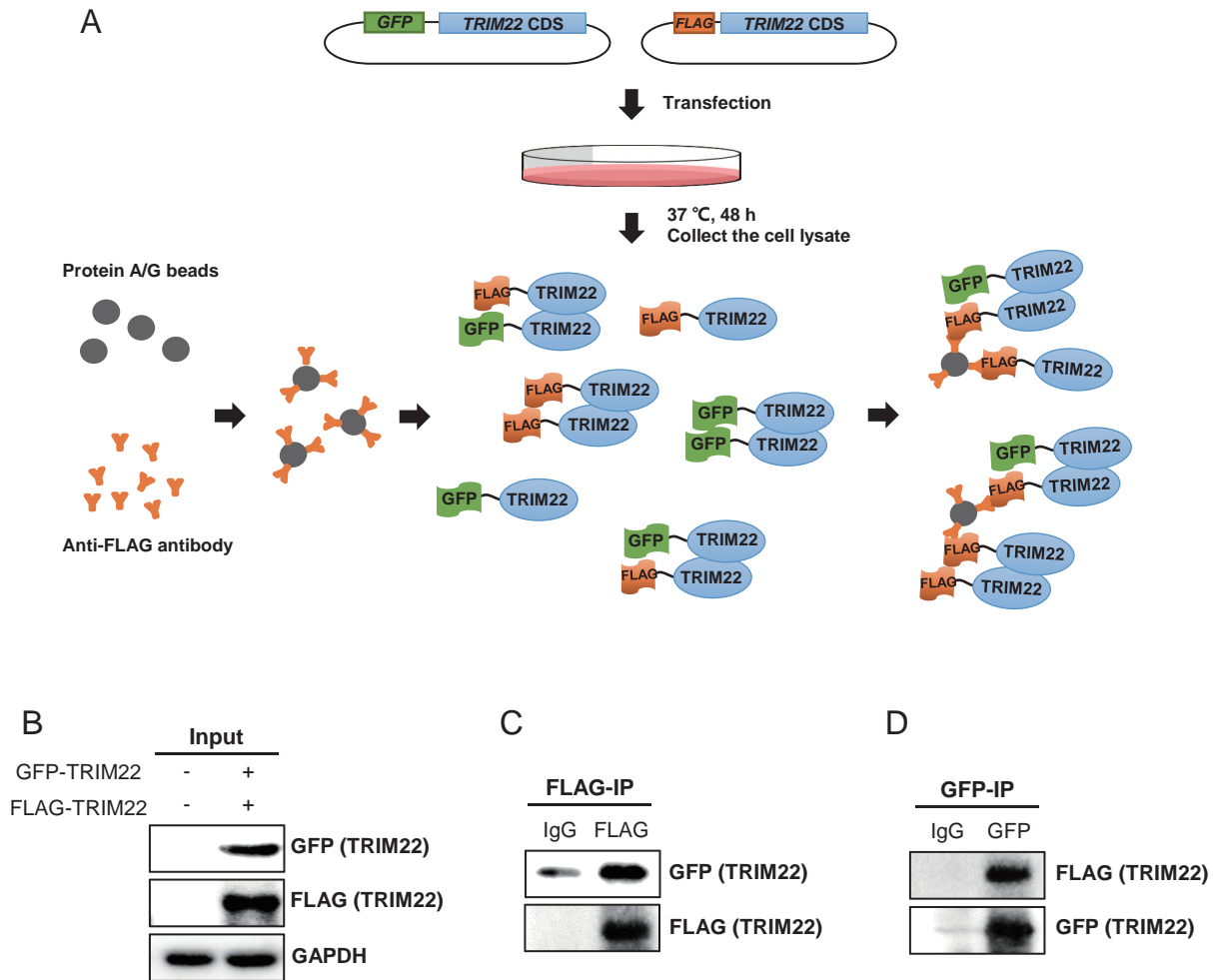
We thus set out to test the binding ability of human TRIM22 using co-immunoprecipitation (Co-IP) (Figure 6A).

First, FLAG-tagged and GFP-tagged human TRIM22 proteins were co-expressed in HEK293T cells (Figure 6B). To prove their interaction, we then pulled down FLAG-tagged TRIM22 from the cell lysate. Western blot showed that when FLAG-tagged TRIM22 was pulled down, GFP-tagged TRIM22 can be detected by anti-GFP antibody (Figure 6C), suggesting that GFP-tagged and FLAG-tagged TRIM22 proteins exist in the same complex in HEK293T cells. To further confirm this binding, we did Co-IP in the opposite way. As expected, FLAG-tagged TRIM22 was also detected in IP of GFP-tagged TRIM22 (Figure 6D). Taken together, our results detected multiple TRIM22 proteins in the same complex, which provides the possibility of its binding with each other.

## Discussion

In this work, we develop QAUST, a method to predict biological functions of protein molecules using three main features: global and local protein structure similarity, PPI information, and functional sequence motifs. In our method, we construct the 3D structure from the amino acid sequence using I-TASSER. Functional analogs are then identified by performing global and local structural similarity search through the functional libraries, with the scoring function involving the confidence score of structural predictions, sequence and structural similarity of the I-TASSER model with the functional templates, and the local active site matches. We have also tried to improve the performance of GO prediction by incorporating PPI information, especially in order to improve the prediction of GO terms under BP and CC aspects. We further developed a novel predictor that extracts functional motifs that are related to a specific GO term and used it as our third predictor.

On a set of 500 non-redundant proteins, QAUST is shown to have higher function prediction accuracy than the other competing methods on most prediction tasks. This performance advantage is mainly a result of combining three different predictors which cover major aspects of proteins. Additionally, our three prediction components complement each other in the sense that they contribute differently to the prediction of the three aspects of GO. While PPI information improves significantly the prediction of BP and CC terms, functional motif detection is mainly useful in improving MF term prediction. However, QAUST has a number of limitations that give room for possible improvement in the future. One main limitation is that QAUST is much more expensive in terms of running time compared to the other methods as reported in “Section 10” in File S1 and Table S4. The second limitation is that our method cannot be directly used to infer functions that are not included in EC or GO systems, since it solely infers



**Figure 6** Experimental validation of homodimerization function of human TRIM22

**A.** Illustration of the Co-IP method to validate the homodimerization of TRIM22. The FLAG-tagged and GFP-tagged human TRIM22 proteins were co-expressed in HEK293T cells by co-transfecting two plasmids into the cells. If TRIM22 forms a homodimer, when FLAG-tagged TRIM22 or GFP-tagged TRIM22 is pulled down, both FLAG-tagged and GFP-tagged TRIM22 should be detected by the corresponding antibodies (4 combinations in total). **B.** Western blot showing the co-expression of FLAG-tagged and GFP-tagged TRIM22 proteins in HEK293T cells. Cells co-transfected with FLAG-tagged and GFP-tagged empty vectors were used as a negative control. **C.** Western blot for FLAG-IP assay. Both FLAG-tagged and GFP-tagged TRIM22 proteins were detected by the corresponding antibodies, whereas mouse IgG was used as a negative control. **D.** Western blot for GFP-IP assay. Both FLAG-tagged and GFP-tagged TRIM22 proteins detected by the corresponding antibodies, whereas mouse IgG was used as a negative control. TRIM22, tripartite motif-containing 22; Co-IP, co-immunoprecipitation.

protein functions from existing protein annotations. Finally, given that the three components we used work differently in predicting different aspects of GO, it may be helpful to weight their scores differently depending on the nature of the GO term evaluated instead of combining the scores in a simple consensus. In particular, advanced machine learning methods, such as deep learning [56–60], could help weight and combine the scores in a more efficient way to obtain better prediction results which could be a possible future improvement of this work.

## Code availability

QAUST can be accessed at <http://www.cbrc.kaust.edu.sa/>

[qaust/submit/](http://qaust/submit/).

## CRedit author statement

**Fatima Zohra Smaili:** Conceptualization, Methodology, Validation, Writing - original draft. **Ambrish Roy:** Writing - original draft, Data curation. **Meshari Alazmi:** Software. **Stefan T. Arold:** Writing - review & editing. **Srayanta Mukherjee:** Writing - original draft. **P. Scott Hefty:** Writing - original draft. **Shuye Tian:** Validation, Writing - original draft. **Wei Chen:** Validation, Supervision. **Xin Gao:** Conceptualization, Methodology, Validation, Writing - review & editing, Supervision. All authors have read and approved the final manuscript.



## Competing interests

The authors have declared no competing interest.

## Acknowledgments

We thank Mr. Chengxin Zhang, Dr. Wei Zhang, and Prof. Yang Zhang for helpful discussions. This work was supported by the King Abdullah University of Science and Technology (KAUST) Office of Sponsored Research (OSR) (Grant Nos. URF/1/1976-04 and URF/1/1976-06).

## Supplementary material

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.gpb.2021.02.001>.

## ORCID

0000-0001-6439-0659 (Fatima Zohra Smaili)

0000-0002-1832-5765 (Shuye Tian)

0000-0003-1738-0591 (Ambrish Roy)

0000-0001-9074-1029 (Meshari Alazmi)

0000-0001-5278-0668 (Stefan T. Arold)

0000-0002-2750-5071 (Srayanta Mukherjee)

0000-0002-2303-2465 (P. Scott Hefty)

0000-0003-3263-1627 (Wei Chen)

0000-0002-7108-3574 (Xin Gao)

## References

- [1] Consortium U. UniProt: a hub for protein information. *Nucleic Acids Research* 2014;43:D204–12.
- [2] Boutet E, Lieberherr D, Tognolli M, Schneider M, Bansal P, Bridge AJ, et al. UniProtKB/Swiss-Prot, the Manually Annotated Section of the UniProt KnowledgeBase: How to Use the Entry View. *Methods Mol Biol* 2016;1374:23–54.
- [3] Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 1997;25:3389–402.
- [4] Eddy SR. Profile hidden markov models. *Bioinformatics* 1998;14:755–63.
- [5] Söding J. Protein homology detection by HMM-HMM comparison. *Bioinformatics* 2005;21:951–60.
- [6] Karplus K, Barrett C, Hughey R. Hidden Markov models for detecting remote protein homologies. *Bioinformatics* 1998;14:846–56.
- [7] Tian W, Skolnick J. How well is enzyme function conserved as a function of pairwise sequence identity? *J Mol Biol* 2003;333:863–82.
- [8] Rost B. Enzyme function less conserved than anticipated. *J Mol Biol* 2002;318:595–608.
- [9] Finn RD, Attwood TK, Babbitt PC, Bateman A, Bork P, Bridge AJ, et al. InterPro in 2017—beyond protein family and domain annotations. *Nucleic Acids Res* 2017;45:D190–9.
- [10] de Lima Morais DA, Fang H, Rackham OJL, Wilson D, Pethica R, Chothia C, et al. SUPERFAMILY 1.75 including a domain-centric gene ontology method. *Nucleic Acids Res* 2011;39:D427–34.
- [11] Rentzsch R, Orengo CA. Protein function prediction using domain families. *BMC Bioinformatics* 2013;14:S5.
- [12] López G, Rojas A, Tress M, Valencia A. Assessment of predictions submitted for the CASP7 function prediction category. *Proteins* 2007;69:165–74.
- [13] Zhang Y. Protein structure prediction: when is it useful? *Curr Opin Struct Biol* 2009;19:145–55.
- [14] Skolnick J, Fetrow JS, Kolinski A. Structural genomics and its importance for gene function analysis. *Nat Biotechnol* 2000;18:283–7.
- [15] Aloy P, Querol E, Aviles FX, Sternberg MJE. Automated structure-based prediction of functional sites in proteins: applications to assessing the validity of inheriting protein function from homology in genome annotation and to protein docking. *J Mol Biol* 2001;311:395–408.
- [16] Roy A, Srinivasan N, Gowri VS. Molecular and structural basis of drift in the functions of closely-related homologous enzyme domains: implications for function annotation based on homology searches and structural genomics. *In Silico Biol* 2009;9:S41–55.
- [17] Bork P, Sander C, Valencia A. Convergent evolution of similar enzymatic function on different protein folds: the hexokinase, ribokinase, and galactokinase families of sugar kinases. *Protein Sci* 1993;2:31–40.
- [18] Spriggs RV, Artymiuk PJ, Willett P. Searching for patterns of amino acids in 3D protein structures. *J Chem Inf Comput Sci* 2003;43:412–21.
- [19] Kinoshita K, Nakamura H. Identification of protein biochemical functions by similarity search using the molecular surface database eF-site. *Protein Sci* 2003;12:1589–95.
- [20] Chang DTH, Chen CY, Chung WC, Oyang YJ, Juan HF, Huang HC. ProteMiner-SSM: a web server for efficient analysis of similar protein tertiary substructures. *Nucleic Acids Res* 2004;32:W76–82.
- [21] Gherardini PF, Helmer-Citterich M. Structure-based function prediction: approaches and applications. *Brief Funct Genomic Proteomic* 2008;7:291–302.
- [22] Laskowski RA, Watson JD, Thornton JM. Protein function prediction using local 3D templates. *J Mol Biol* 2005;351:614–26.
- [23] Zhang C, Zheng W, Freddolino PL, Zhang Y. MetaGO: predicting gene ontology of non-homologous proteins through low-resolution protein structure prediction and protein–protein network mapping. *J Mol Biol* 2018;430:2256–65.
- [24] Roy A, Yang J, Zhang Y. COFACTOR: an accurate comparative algorithm for structure-based protein function annotation. *Nucleic Acids Res* 2012;40:W471–7.
- [25] Brun C, Chevenet F, Martin D, Wojcik J, Guénoche A, Jacq B. Functional classification of proteins for the prediction of cellular function from a protein-protein interaction network. *Genome Biol* 2003;5:R6.
- [26] Chua HN, Sung WK, Wong L. Using indirect protein interactions for the prediction of Gene Ontology functions. *BMC Bioinformatics* 2007;8:S8.
- [27] Barabási AL, Gulbahce N, Loscalzo J. Network medicine: a network-based approach to human disease. *Nat Rev Genet* 2011;12:56–68.
- [28] Magnez R, Thiroux B, Taront S, Segaula Z, Quesnel B, Thuru X. PD-1/PD-L1 binding studies using microscale thermophoresis. *Sci Rep* 2017;7:17623.
- [29] Lan L, Djuric N, Guo Y, Vucetic S. MS-k NN: protein function prediction by integrating multiple data sources. *BMC Bioinformatics* 2013;14:S8.
- [30] You R, Zhang Z, Xiong Y, Sun F, Mamitsuka H, Zhu S. GOLabeler: improving sequence-based large-scale protein function prediction by learning to rank. *Bioinformatics* 2018;34:2465–73.
- [31] Kulmanov M, Khan MA, Hoehndorf R, Wren J. DeepGO: predicting protein functions from sequence and interactions using a deep ontology-aware classifier. *Bioinformatics* 2017;34:660–8.
- [32] Gong Q, Ning W, Tian W. GoFDR: a sequence alignment based method for predicting protein functions. *Methods* 2016;93:3–14.

- [33] Zou Z, Tian S, Gao X, Li Y. mlDEEPre: multi-functional enzyme function prediction with hierarchical multi-label deep learning. *Front Genet* 2019;9:714.
- [34] Gao X, Bu D, Xu J, Li M. Improving consensus contact prediction via server correlation reduction. *BMC Struct Biol* 2009;9:28.
- [35] Chen P, Hu SS, Zhang J, Gao X, Li J, Xia J, et al. A sequence-based dynamic ensemble learning system for protein ligand-binding site prediction. *IEEE-ACM Transactions on Computational Biology and Bioinformatics* 2016;13:901–12.
- [36] Chen P, Huang JZ, Gao X. LigandRFs: random forest ensemble to identify ligand-binding residues from sequence information alone. *BMC Bioinformatics* 2014;15:S4.
- [37] Chen P, Li J, Wong L, Kuwahara H, Huang JZ, Gao X. Accurate prediction of hot spot residues through physicochemical characteristics of amino acid sequences. *Proteins* 2013;81:1351–62.
- [38] Zhang Y. I-TASSER server for protein 3D structure prediction. *BMC Bioinformatics* 2008;9:40.
- [39] Zhang C, Freddolino PL, Zhang Y. COFACTOR: improved protein function prediction by combining structure, sequence and protein–protein interaction information. *Nucleic Acids Res* 2017;45:W291–9.
- [40] Szklarczyk D, Franceschini A, Wyder S, Forslund K, Heller D, Huerta-Cepas J, et al. STRING v10: protein–protein interaction networks, integrated over the tree of life. *Nucleic Acids Res* 2015;43:D447–52.
- [41] Webb EC. Recommendations of the Nomenclature Committee of the International Union of Biochemistry and Molecular Biology on the Nomenclature and Classification of Enzymes. Pittsburgh: Academic Press;1992.
- [42] Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene Ontology: tool for the unification of biology. *Nat Genet* 2000;25:25–9.
- [43] Wu S, Zhang Y. LOMETS: a local meta-threading-server for protein structure prediction. *Nucleic Acids Res* 2007;35:3375–82.
- [44] Jiang Y, Oron TR, Clark WT, Bankapur AR, D’Andrea D, Lepore R, et al. An expanded evaluation of protein function prediction methods shows an improvement in accuracy. *Genome Biol* 2016;17:17:184.
- [45] Piovesan D, Giollo M, Leonardi E, Ferrari C, Tosatto SCE. INGA: protein function prediction combining interaction networks, domain assignments and sequence similarity. *Nucleic Acids Res* 2015;43:W134–40.
- [46] Porter CT, Bartlett GJ, Thornton JM. The Catalytic Site Atlas: a resource of catalytic sites and residues identified in enzymes using structural data. *Nucleic Acids Res* 2004;32:D129–133.
- [47] Zhang Y, Skolnick J. TM-align: a protein structure alignment algorithm based on the TM-score. *Nucleic Acids Res* 2005;33:2302–9.
- [48] UniProt Consortium. The Universal Protein Resource (UniProt) in 2010. *Nucleic Acids Res* 2010;38:D142–8.
- [49] Li Y, Wang S, Umarov R, Xie B, Fan M, Li L, et al. DEEPre: sequence-based enzyme EC number prediction by deep learning. *Bioinformatics* 2017;34:760–9.
- [50] Radivojac P, Clark WT, Oron TR, Schnoes AM, Wittkop T, Sokolov A, et al. A large-scale evaluation of computational protein function prediction. *Nat Methods* 2013;10:221–7.
- [51] Mueser TC, Jones CE, Nossal NG, Hyde CC. Bacteriophage T4 gene 59 helicase assembly protein binds replication fork DNA. The 1.45 Å resolution crystal structure reveals a novel  $\alpha$ -helical two-domain fold. *J Mol Biol* 2000;296:597–612.
- [52] Barr SD, Smiley JR, Bushman FD. The interferon response inhibits HIV particle production by induction of TRIM22. *PLoS Pathog* 2008;4:e1000007.
- [53] Di Pietro A, Kajaste-Rudnitski A, Oteiza A, Nicora L, Towers GJ, Mehti N, et al. TRIM22 inhibits influenza A virus infection by targeting the viral nucleoprotein for degradation. *J Virol* 2013;87:4523–33.
- [54] Yang C, Zhao X, Sun D, Yang L, Chong C, Pan Y, et al. Interferon alpha (IFN $\alpha$ )-induced TRIM22 interrupts HCV replication by ubiquitinating NS5A. *Cell Mol Immunol* 2016;13:94–102.
- [55] Lou J, Wang Y, Zheng X, Qiu W. TRIM22 regulates macrophage autophagy and enhances *Mycobacterium tuberculosis* clearance by targeting the nuclear factor–multiplicity  $\kappa$ B/beclin 1 pathway. *J Cell Biochem* 2018;119:8971–80.
- [56] Xia Z, Li Y, Zhang B, Li Z, Hu Y, Chen W, et al. DeeReCT-PolyA: a robust and generic deep learning method for PAS identification. *Bioinformatics* 2019;35:2371–9.
- [57] Umarov R, Kuwahara H, Li Y, Gao X, Solovyev V. Promoter analysis and prediction in the human genome using sequence-based deep learning models. *Bioinformatics* 2019;35:2730–7.
- [58] Kim JS, Gao X, Rzhetsky A. RIDDLE: race and ethnicity imputation from disease history with deep learning. *PLoS Comput Biol* 2018;14:e1006106.
- [59] Li Y, Xu F, Zhang F, Xu P, Zhang M, Fan M, et al. Dlbi: deep learning guided bayesian inference for structure reconstruction of super-resolution fluorescence microscopy. *Bioinformatics* 2018;34:i284–94.
- [60] Li Y, Huang C, Ding L, Li Z, Pan Y, Gao X. Deep learning in bioinformatics: introduction, application, and perspective in the big data era. *Methods* 2019;166:4–21.