

1 **KAUST: Protein Function Prediction Using Structure Similarity**
2 **Search, Protein Interaction and Functional Sequence Motifs**

3 Fatima Zohra Smaili^{1,a,#}, Shuye Tian^{2,b,#}, Ambrish Roy^{3,c}, Meshari Alazmi^{1,4,d}, Stefan T.
4 Arold^{5,e}, Srayanta Mukherjee^{3,f}, P. Scott Hefty^{6,g}, Wei Chen^{2,h,*}, Xin Gao^{1,i,*}

5
6 ¹*King Abdullah University of Science and Technology (KAUST), Computational*
7 *Bioscience Research Center (CBRC), Computer, Electrical and Mathematical Sciences*
8 *and Engineering (CEMSE) Division, Thuwal 23955, Saudi Arabia*

9 ²*Department of Biology, Southern University of Science and Technology of China*
10 *(SUSTC), Shenzhen 518055, China*

11 ³*Department of Computational Medicine and Bioinformatics, University of Michigan,*
12 *Ann Arbor, MI 48109, USA*

13 ⁴*College of Computer Science and Engineering, University of Hail, Hail, Saudi Arabia*

14 ⁵*King Abdullah University of Science and Technology (KAUST), Biological and*
15 *Environmental Sciences and Engineering (BESE) Division, Thuwal 23955, Saudi Arabia.*

16 ⁶*Department of Molecular Bioscience, University of Kansas, 2030 Becker Dr, Lawrence,*
17 *KS 66047, USA*

18

19 #Equal contribution.

20 *Corresponding authors.

21 Email: chenw@sustech.edu.cn (Chen W), xin.gao@kaust.edu.sa (Gao X).

22

23 **Running title:** *Smaili FZ et al / Protein Function Prediction*

24

25 ^a ORCID: 0000-0001-6439-0659.

26 ^b ORCID: 0000-0002-1832-5765.

27 ^c ORCID: 0000-0003-1738-0591.

28 ^d ORCID: 0000-0001-9074-1029.

29 ^e ORCID: 0000-0001-5278-0668.

30 ^f ORCID: 0000-0002-2750-507.

31 ^g ORCID: 0000-0002-2303-2465.

32 ^h ORCID: 0000-0003-3263-1627.

33 ⁱORCID: 0000-0002-7108-3574.

34

35

36 This manuscript contains a total of 8739 words, 6 figures, 3 tables, 4 supplementary

37 figures and 4 supplementary tables.

38 **Abstract**

39 The number of available protein sequences in public databases is increasing
40 exponentially. However, a significant percentage of these sequences lack functional
41 annotation, which is essential for the understanding of how biological systems operate.

42 We propose a novel method, QAUST, to infer protein functions, specifically Gene
43 Ontology (GO) terms and Enzyme Commission (EC) numbers. Our method uses three
44 sources of information: structure information encoded by global and local structure
45 similarity search, biological network information inferred by protein-protein interaction
46 data, and sequence information extracted from functionally discriminative sequence
47 motifs. The three pieces of information are combined by consensus averaging to make the
48 final prediction. Our approach has been tested on 500 protein targets from the CAFA
49 benchmark set. The results show that our method provides accurate functional annotation
50 and outperforms other prediction methods based on sequence similarity search or
51 threading. We further demonstrate that a previously unknown function of TRIM22
52 protein predicted by QAUST can be experimentally validated. Availability:
53 <http://www.cbrc.kaust.edu.sa/qaust/submit/>.

54

55 **KEYWORDS:** Protein function prediction; Gene Ontology (GO) terms; Enzyme
56 Commission (EC) numbers; Protein structure similarity; Protein-protein interaction;
57 Functionally discriminative motifs

58

59 **Introduction**

60 As of today, over 88 million protein sequences are available in the UniProtKB/TrEMBL [1]
61 database. However, this increase in the number of known protein sequences does not
62 reflect a parallel increase in our biological knowledge, as less than 1% of these sequences
63 have a manually annotated function [2]. On the other hand, the functional annotation of
64 these sequences is not only an essential step for the understanding of physiological
65 processes and biological systems in living entities, but also one of the highly challenging
66 tasks in biology, which is why there is an increasing need to provide reliable, automated
67 protein function annotation.

68 Significant efforts have been made to identify evolutionarily related proteins and
69 automatically transfer functional annotations between homologous protein pairs [3–6].
70 To make such sequence similarity based functional transfer possible, powerful sequence-
71 alignment methodologies have been developed. In particular, algorithms like
72 BLAST/PSI-BLAST [3] and hidden Markov model (HMM) based techniques [4–6] have
73 been frequently used to transfer functional annotations between homologous proteins.
74 The underlying assumption of these sequence-based methods is that evolutionarily related
75 proteins may inherit the function of a shared common ancestor. However, there are
76 numerous cases in which proteins with high sequence similarity have distinct functions
77 [7,8]. To partially address the problem, several methods have been developed to predict
78 function using annotated conserved sequence motifs that are responsible for the
79 functional aspect of the protein. These methods typically construct the sequence motifs
80 from multiple sequence alignment of proteins belonging to the same protein family with
81 known function [9–11]. They, however, have two major limitations. First, high-quality
82 sequence alignment is typically required for motif construction, which is not trivial to
83 obtain especially when the sequence homology is low. Second, the accuracy is limited by
84 the quality of functional annotation of motifs. To overcome these limitations, we propose
85 in this work to use a protein-specific “functionally discriminative motif” constructed
86 from sequence fragments excised from the template sequence.

87 From another perspective, the 3D structure of a protein sequence is believed to be
88 more involved in its biological function [12,13] since structures are more conserved than
89 the sequences. The 3D structure of a protein can therefore provide additional information

90 for function transfer, especially when the sequence similarity between related proteins is
91 too low for sequence homolog detection [14,15]. However, the relationship between the
92 protein function and its structure is not straightforward, as in some cases, similar
93 structures perform the same function while in many cases similar folds perform different
94 functions [16,17]. Therefore, many prediction methods have been relying on local
95 structure similarity search methods rather than global similarity search to identify
96 functionally homologous proteins [18–20]. Most of these approaches scan the query
97 protein against a library of known conserved spatial motifs or known active sites (*e.g.*,
98 binding sites) with known function [21]. Local similarity search methods have been
99 proven to be quite accurate in detecting functional similarity between proteins of different
100 folds, but they also have a high probability of producing false positive matches [22]. One
101 possible solution is to combine global and local structure alignment to overcome the
102 promiscuity of global structure comparison and low specificity of local structure
103 matching [23,24], which we implement in this project.

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

111 Finally, recently there is an emergence of methods which combine multiple sources
112 of information (PPI, domains, sequence alignments, *etc.*) using advanced machine
113 learning algorithms to perform function prediction. These methods have shown to
114 improve the prediction performance over methods that use only one type of information
115 [29–37].

116 In this work, we propose a new protein function prediction method, Quantitative
117 Annotation of Unknown STructure (QAUST), which combines the global and local
118 structure similarity search with protein-protein interaction networks and functional
119 sequence motif detection. Our approach follows a sequence-to-structure-to-function
120 workflow. Starting from the protein amino acid sequence, we first generate structure

121 predictions by the Iterative Threading ASSEMBLY Refinement method (I-TASSER) [38].
122 The predicted structure is then used to identify the proteins with similar functions based
123 on a combination of global and local structure similarity search method that follows the
124 same pipeline used in COFACTOR [24,39]. Protein-protein interaction information is
125 meanwhile extracted from the STRING database [40]. And finally, we extract
126 functionally discriminative sequence motifs as our third main prediction feature. The
127 confidence scores obtained from these three features are combined in a consensus
128 function to obtain our final confidence score.

129 Since the terminology of a “protein function” might be ambiguous, we would like to
130 clarify that the definitions of function followed in this work is Enzyme Commission (EC)
131 numbers [41] and Gene Ontology (GO) terms [42]. EC numbers are used to categorize
132 enzymes into hierarchical families using a numerical classification. Specifically, the EC
133 number (which is composed of four numbers separated by periods *i.e.*, A.B.C.D) refers to
134 the reaction catalyzed by a specific enzyme. On the other hand, the GO terms are a set of
135 controlled vocabulary to formally describe proteins and RNAs based upon their functions.
136 Three aspects of ontologies, Biological Process (BP), Cellular Component (CC) and
137 Molecular Function (MF), are defined in this database. Each one of these three GO
138 aspects is represented by a structured directed acyclic graph (DAG), where nodes
139 represent GO terms which describe gene product functions, while the edges represent the
140 relationships (“is_a” or “part_of”) between the GO terms. In GO’s functional hierarchy,
141 the more general functions are on the top of the graph while more specific terms are
142 usually present further down the graph.

143 Our prediction results are compared to the following programs: COFACTOR
144 [39,43], a global and local structure similarity-based method, LOMETS [44], a meta-
145 threading algorithm, HHsearch [5], an HMM-based method that is widely used to detect
146 protein homologs, BLAST [3], which transfers annotations based on sequence similarity,
147 naïve baseline which predicts GO terms based on their annotation frequency, as well as
148 two highly-ranked methods from the CAFA assessment [45], GoFDR [32], and INGA
149 [46].

150

151 **Methods**

152 **Dataset**

153 To evaluate QAUST for EC prediction, we used the benchmark data set of COFACTOR
154 [39,43] as our testing data set. This data set consists of 318 enzymes with unique EC
155 numbers (first three digits) covering all 6 enzyme classes. Similarly, all sequences in our
156 template libraries with a sequence identity > 30% with the query enzymes are excluded
157 from the template libraries.

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

164

165 **Enzyme Commission (EC) number prediction**

166 *Global and local similarity search*

167 The first step of our protein function prediction is the generation of the predicted 3D
168 model of the query protein using I-TASSER [38] as outlined in Section S1. The predicted
169 model of the query protein obtained from I-TASSER is then scanned against a non-
170 redundant (pairwise sequence identity no more than 90%) structure template library of
171 2,385 enzymes with at least the first three digits of EC number annotated by the Catalytic
172 Site Atlas (CSA) database [47]. This library scanning detects homologous structure
173 templates to the query proteins using two types of structure similarity search programs:
174 *global similarity search* and *local similarity search*.

175

176 *Global similarity search*

177 Templates with a similar global structure to the predicted structure of the query protein
178 are detected from the template library using TM-align [48]. Another important
179 consideration when searching for templates with similar global folds to the query protein
180 is the quality of the structural models. Appraising the accuracy of the structure modeling
181 in the scoring scheme helps to reduce the number of false positive predictions. In this

182 particular case, the quality of the predicted I-TASSER model generated in the previous
183 step is evaluated using *Cscore* [38].

184

185 Local similarity search

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

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

211 In the third step, the best alignment of the local active site residues in the spheres
212 between the query and the template is identified using a scoring function similar to TM-

213 align. Starting from the initial superposition of the query and template protein structures,
214 we perform a Needleman-Wunsch dynamic programming to generate the best alignment
215 for the residues in the selected sphere of the template and the query, where the alignment
216 score matrix S_{ij} for aligning the i th residue in the query and the j th residue in the template
217 is defined as:

$$218 \quad S_{ij} = \left[\frac{1}{1 + \frac{d_{ij}}{d_0}} + M_{ij} \right], \quad (1)$$

219 where d_{ij} is the Ca distance between residues i and j , d_0 is the distance cutoff given by
220 $d_0 = 1.24\sqrt[3]{L-15} - 1.8$ obtained from TM-align, M_{ij} is the substitution score between the
221 i th and j th residues taken from the BLOSUM62 mutation matrix with the value
222 normalized by the diagonal element in the mutation matrix. The gap penalty is set as -1.
223 For a given scoring matrix S_{ij} , a new alignment is generated by dynamic programming. A
224 new superposition and scoring matrix are then constructed based on the new alignment to
225 obtain a newer alignment from dynamic programming. This procedure is iteratively
226 repeated until the final alignment is converged. For each alignment, the active site match
227 (AcM) is evaluated using an alignment score defined as:

$$228 \quad \text{AcM} = \frac{1}{N_t} \sum_{i=1}^{N_{ali}} \frac{1}{1 + \left(\frac{d_{ii}}{d_0}\right)^2} + \frac{1}{N_t} \sum_{i=1}^{N_{ali}} M_{ii}, \quad (2)$$

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

238

239 Scoring function for global and local similarity search

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

243

$$244 \quad QAUSTEC = C_{norm} \cdot \left[TM + \frac{Cov}{1 + RMSD_{ali}} \right] + 2 \cdot ID_{ali} \cdot Cov + \frac{AcM}{2}. \quad (3)$$

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

$$250 \quad QAUSTEC_{norm} = \frac{2}{1 + \exp(-QAUSTEC)} - 1. \quad (4)$$

251

252 ***Gene Ontology (GO) term prediction***

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

260

261 *Global protein structure similarity*

262 Similar to EC prediction, I-TASSER is also used here to construct the corresponding 3D
263 model to the query sequence. The model obtained is then scanned against a library of
264 templates to identify those which share a similar global structure to the query model
265 (<https://zhanglab.ccmb.med.umich.edu/BioLiP/library.html>). For the time being, the
266 functionally important residues for most of the proteins in the GO template library are
267 unknown. Therefore, only the global similarity search is taken into consideration when
268 sorting the hits from the GO library. Global similarity search for GO prediction is done in

269 a similar way to global similarity search for EC prediction described in the previous
270 section. The only difference is that to select the best hits for GO prediction, we rank a
271 template using the Fh-score defined as:

272

$$273 \quad Fh_{score} = C_{norm} \times \left(TM_{score} + \frac{1}{1 + RMSD_{ali}} \times Cov \right) + 3 \times ID_{ali} \times Cov. \quad (5)$$

274

275 Since each single protein can be annotated with multiple GO terms and the global search
276 may result in many close template structures, a query protein can have multiple GO term
277 predictions with high Fh-scores. Therefore, the confidence score of each GO term is
278 calculated as follow:

$$279 \quad P_{structure}(\lambda) = \frac{1}{N} \sum_{i=1}^{N_\lambda} Fh(i), \quad (6)$$

280 where λ represents a given GO term, N_λ is the number of templates annotated with the
281 GO term λ , and N is the total number of templates selected for generating the consensus.
282 When multiple close templates are available, we only consider the templates with an Fh-
283 score >1 . For those query proteins with less than 10 templates of Fh-score >1 , the top 10
284 templates are selected for generating the consensus prediction regardless of the Fh-score.
285 Also, given the hierarchical nature of the GO DAG, we consider that when a protein is
286 annotated with a given GO term, all its ancestor GO terms (through “is_a” relation) are
287 automatically implied. Therefore, once a GO term λ is scored, we score all its ancestor
288 terms as well. The score of any ancestor GO term μ of term λ is calculated as:

$$289 \quad P_{structure}(\mu) = P_{structure}(\lambda) \times \left(1 + \frac{N_\mu}{N_0} \right), \quad (7)$$

290 where N_μ and N_0 are the number of leaf nodes under node μ and the root node,
291 respectively. Since COFACTOR [39,43] uses a similar structure scoring function, we
292 have highlighted the main differences between QAUST and COFACTOR in Section S2,
293 and compared our method with COFACTOR in the experiments.

294

295 *Protein-protein interaction network*

296 We exploit the information provided by the STRING [40] database, which is a library of
297 PPI networks, to extend our prediction set. The query protein sequence is mapped to its
298 corresponding STRING entry by BLAST, with minimum sequence identity cutoff of 90%.
299 Extracting the PPI partners of the query, we calculate the confidence score of STRING
300 for a GO term λ ($P_{STRING}(\lambda)$) as the frequency of the GO term λ among the
301 experimentally annotated interaction partners of the query protein:

$$302 \quad P_{STRING}(\lambda) = \frac{n_{\lambda}}{N}, \quad (8)$$

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

307

308 *Functionally discriminative sequence motifs*

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

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

317 Detection of sequence templates for query protein: The sequence homologs of the
318 query sequence are detected by PSI-BLAST [3] from the Uniref90 database [49]. We
319 filter all obtained homologs with sequence identity $> 30\%$ to the query.

320 Identification of functionally discriminative motifs given a GO term: We map all the
321 selected sequence homologs of the query to their corresponding GO annotations in the
322 UniProt-GOA database (<http://www.ebi.ac.uk/GOA>). GO terms assigned with “Inferred
323 from Electronic Annotation” (IEA) or “No biological Data available” (ND) evidence
324 codes are not considered. We also filter out annotations with evidence code IPI (Inferred
325 from Physical Interactions) since we use protein-protein interaction (PPI) information in
326 our features. After filtering these annotations, we are left with the annotations based on

327 evidence codes: EXP, IDA, IMP, IGI, IEP, TAS, and IC. For each GO term λ , we build
328 two sets of sequences from the set of homolog sequences detected in the previous step.
329 These two sets are: the “annotated set”, which is the set of sequence homologs annotated
330 with this specific GO term, and the “not-annotated set”, which is the set of sequence
331 homologs not annotated with this given GO term. For each one of these two sets, we
332 extract the ten most frequent motifs by extracting all unique amino acid motifs of length
333 [4,7] from the sequence set using sliding windows. These motifs are ranked in
334 descending order by their occurrences. The top 10 most frequent motifs are the initial
335 "frequent list", while the remaining motifs are in the "waiting list". If, within the
336 "frequent list", a short motif is a substring of another longer motif, the shorter motif is
337 discarded, and the most frequent motif from the "waiting list" is transferred to "frequent
338 list" to ensure that the latter always has 10 motifs. This process is iterated until, in the
339 "frequent list", any motif is not a substring of another motif. The motifs in the "frequent
340 list" are used for matching the query in the next step.

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

$$345 \quad P_{MOTIF}(\lambda) = \frac{n_q(\lambda)}{N(\lambda)} \left[1 - \frac{n_q(\lambda^c)}{N(\lambda^c)} \right], \quad (9)$$

346 where λ is the given GO term, $N(\lambda)$ and $N(\lambda^c)$ are the number of frequent patterns from
347 the “annotated set” and “not-annotated set”, both of which equal to 10. $n(\lambda)$ and $n(\lambda^c)$
348 are the corresponding number of matched patterns at the query sequence. This score can
349 take any value from 0 to 1. An ideal value of this score would be equal to 1, which
350 happens when all the sequences in the annotated set contain these frequent motifs and
351 none of the sequences in the not-annotated set contains these same motifs. This scoring
352 function has been designed to penalize the prediction in case the query sequence matches
353 a high number of frequent motifs from the not-annotated set. This way, the scoring
354 function accounts for two essential pieces of information: which set has the maximum
355 number of frequent motifs matched in the query, and how significant is the difference
356 between the number of matched motifs from the annotated set to that from the not-

357 annotated set. **Figure 2** shows a flowchart detailing the three steps of extracting
358 functional sequence motifs.

359

360 Consensus

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

364

$$365 \quad P_{consensus}(\lambda) = 1 - \prod_{m \in \{structure, STRING, MOTIF\}} (1 - P_m(\lambda)) \quad (10)$$

366 This equation used to calculate the consensus has been previously used by other methods
367 for protein function prediction [46]. If one or more predictors are not available for a given
368 term (*e.g.*, no interaction partners are known for the given query), only the available
369 predictors are used to obtain the confidence score. Also, since GO uses the true-path rule
370 (*i.e.*, if a protein is associated by a term, it is also implicitly annotated by its ancestors),
371 for every predicted GO term, all its ancestors are considered to be predicted as well since
372 they are more general terms.

373

374 **Results**

375 **Prediction of EC numbers**

376 We compared the EC prediction performance of our method to five methods: HHsearch
377 [5], LOMETS [44], BLAST [3], COFACTOR [39,43] and DEEPre webserver [50]. We
378 compared the performance of these methods based on precision (positive predictive value)
379 and recall (sensitivity) rates. **Figure 3** shows the precision-recall graph corresponding to
380 four baseline methods as well as QAUST. Since the DEEPre webserver does not report
381 the confidence score with the annotation, we could not draw the precision-recall curves
382 but compared QAUST to DEEPre based on accuracy. An EC number prediction is
383 considered to be “true” if the first three digits of the EC number from the hit are identical
384 to those of the query protein; otherwise the hit is considered to be “false”. As shown in
385 Figure 3, the rate of true positive predictions using the EC-score is much higher than that
386 of HHsearch, LOMETS, BLAST and COFACTOR at most recall rates. QAUST has also

387 an area under precision-recall curve (AUPRC) of 0.712 which is higher than that of
388 COFACTOR (0.643), LOMETS (0.510), and HHsearch (0.489). **Table 1** reports the
389 accuracy of QAUST compared to five other methods including DEEPre and
390 COFACTOR. The results show that DEEPre has a slightly higher performance than
391 QAUST in terms of accuracy, which is probably due to the fact that DEEPre is a machine
392 learning method trained on a large number of enzymes with known functions that overlap
393 or contain close homologs to our test data.

394

395 **Prediction of GO terms**

396 To assess the contribution of individual predictors to the GO prediction performance by
397 QAUST, we visualize the precision-recall curve of the structure similarity search alone
398 ($P_{\text{structure}}$), the precision-recall curve of structure similarity search combined with PPI
399 information ($P_{\text{structure}}$ and P_{STRING}), and that of the final QAUST prediction ($P_{\text{structure}}$,
400 P_{STRING} and P_{MOTIF}). Additionally, we compared the prediction performance on our
401 dataset (please see subsection Dataset under section Methods) to COFACTOR [39,43],
402 BLAST [3], LOMETS [44], HHsearch [5], INGA [46] webserver, a method that
403 combines BLAST, PPI information and Pfam in one predictor, and GoFDR [32], one of
404 the top function prediction methods at the CAFA assessment [45] which uses a machine
405 learning model as classifier and discriminative residues as the main feature.

406 The performance has been primarily evaluated using precision-recall curves
407 computed at each prediction score threshold. We also used the F_{max} measure as a
408 quantitative measure to evaluate the overall performance of the precision-recall curves.
409 Precision, recall, and F_{max} are defined in the same way as the CAFA evaluation [51]. The
410 F_{max} measure has been computed as the maximum value of the F_{measure} which is computed
411 at each threshold as $\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$.

412 Precision at threshold t is defined as $\frac{|P_x(t) \cap C_x|}{|P_x(t)|}$, while recall is defined as $\frac{|P_x(t) \cap C_x|}{|C_x|}$
413 where x is a query protein, $P_x(t)$ is the set of predicted terms for x at threshold t and C_x is
414 the set of correct terms that x is experimentally annotated with.

415 Similar to the CAFA evaluation [34], we also reported the minimum semantic
416 distance (S_{min}) as an additional evaluation metric for GO prediction. S_{min} is defined as

417 $\min_t \{\sqrt{ru(t)^2 + mi(t)^2}\}$, where $ru(t)$ is the remaining uncertainty at threshold t defined
418 as $\frac{1}{n_e} \sum_{i=1}^{n_e} \sum_f ic(f) \mid f \notin P_i(t) \wedge f \in C_i$ and $mi(t)$ is the misinformation at threshold t
419 defined as $\frac{1}{n_e} \sum_{i=1}^{n_e} \sum_f ic(f) \mid f \in P_i(t) \wedge f \notin C_i$, where n_e is the number of proteins in
420 our dataset, $P_i(t)$ is the set of predicted GO term for protein i at threshold t , C_i is the set
421 of terms that protein i is actually annotated with and $ic(f)$ is the information content of
422 the GO term f . The very general and unspecific GO terms such as “Molecular Function”,
423 “Biological Process”, “Cellular Component”, “Binding” and “Protein Binding” are
424 excluded from the evaluation.

425 As shown in **Figure 4**, our method combining structure, PPI and functional motif
426 information achieves higher precision than most other methods at most recall points, in
427 particular for MF and BP. For our dataset, structure and motif information has been used
428 for all proteins. However, the PPI information from STRING was missing for 74 proteins.
429 In this case, only the structure and motif information is used. The F_{\max} measure values are
430 reported in **Table 2** while the minimum semantic distance (S_{\min}) values are reported in
431 **Table 3**. Surprisingly, in CC prediction, naïve baseline which predicts GO terms based
432 on their annotation frequency, achieves higher performance than all other methods
433 including QAUST. In fact, in the CAFA assessment [45], naïve baseline also
434 outperformed most of the other methods in predicting CC terms. One possible
435 explanation for why naïve baseline has higher F_{\max} for CC terms prediction is because the
436 most frequently used CC terms in protein annotation are usually part of a small set of
437 very general terms such as “cytoplasm” or “intracellular part”. Since the naïve baseline is
438 solely based on frequency, it increases the chance of predicting a true positive [45]. We
439 also reported the p-values obtained from the Mann-Whitney U test to assess the
440 significance of the difference in performance of QAUST compared to all other methods
441 in Section S3 and Table S1.

442 As a further analysis, to investigate if the performance of our method is solely due to
443 the power of the I-TASSER structure prediction we used, we replaced the I-TASSER
444 structure prediction component of our method by HHsearch and LOMETS structure
445 prediction, respectively. Our results show that no matter which structure prediction
446 method is used, our scoring function, $P_{\text{structure}}$, can significantly improve the performance

447 on predicting the GO terms. Meanwhile, among the three structure prediction methods, I-
448 TASSER with $P_{\text{Structure}}$ consistently performs the best over all three GO hierarchy
449 branches of MF, BP and CC, whereas LOMETS with $P_{\text{Structure}}$ has the second best
450 performance on MF and CC, and HHsearch with $P_{\text{Structure}}$ is the second best on predicting
451 BP terms (Section S4 and Figure S1). Additionally, we have evaluated the performance
452 of our method when only PPI and motif information are used without including any
453 structure-based information. The results show that the function prediction performance
454 drops when structure features are not used (Section S5 and Figure S2).

455

456 *How do protein-protein interaction information and functional sequence motifs improve*
457 *the prediction?*

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

469 In addition to the structure similarity search and the protein-protein interaction
470 features, the results show that the functional motifs extracted improve the performance of
471 the prediction significantly. As a further analysis, we have evaluated the performance of
472 our functional motif detection method when both predicted and experimentally annotated
473 GO terms are taken into consideration instead of considering experimental annotations
474 only. The results of this experiment are reported in Section S6 and Table S2. In addition
475 to comparing the performance of our method to BLAST, LOMETS, HHsearch, and
476 COFACTOR, we also compared it to INGA and GoFDR, two top methods from CAFA

477 [45] in particular for MF and BP terms prediction, and to naïve baseline which is one of
478 the performance references used in CAFA.

479

480 **Case study**

481 To better illustrate the performance of QAUST and the contribution of each component
482 to the prediction, we used as an example Bacteriophage T4 gene 59 helicase assembly
483 protein (P13342) (the cyan structure in **Figure 5A**), which is a DNA binding protein
484 required mainly for DNA replication in the late stage of T4 infection [52]. Figure 5B
485 shows the set of BP terms associated with this protein. In this particular example, both
486 BLAST and INGA did not predict any correct term for this protein (the naïve root term is
487 not counted here). When solely using global structure similarity ($P_{\text{structure}}$), we could only
488 predict one single correct BP term. This makes sense because all the queries in our test
489 set are difficult targets, which do not have close homologs in the template database. For
490 instance, the closest template for this query P13342 is the methionine-tRNA ligase (the
491 magenta structure in Figure 5A), which corresponds to the PDB ID 2CT8A. The
492 sequence identity between P13342 and 2CT8A is only 6.84% and the TM-score between
493 the two structures is only 0.24. Therefore, structure similarity or homology-based
494 methods are not expected to predict the function of the query well. Structure information
495 ($P_{\text{Structure}}$) combined with PPI predicted three correct terms out of six. On the other hand,
496 QAUST predicted four correct terms out of six. In addition, the prediction of QAUST is
497 at least one level deeper in the GO hierarchy than the other methods. Meanwhile, the
498 predicted MF and CC terms for this protein by QAUST are at least as accurate as other
499 methods. Two other case studies illustrating examples for predicting MF and CC terms
500 can be found in Sections S7 and S8, respectively.

501

502 **Experimental validation of TRIM22 dimerization**

503 To provide an experimental assessment of the performance of QAUST, we chose the
504 human tripartite motif-containing 22 (TRIM22) protein as an example. TRIM22 is known
505 as an interferon-inducible protein which shows antiviral activity, such as HIV, HBV and
506 HCV [53–55]. Recent studies also showed that TRIM22 mediates autophagy in human

507 macrophages [56]. However, the function of TRIM22 is still not comprehensively
508 understood as the protein only exists in primates.

509 We applied QAUST to predict the function for TRIM22. Among the predicted GO
510 terms with high consensus scores (Section S9 and Table S3), some of the CC and BP
511 terms agree well with the previously known functions of TRIM22, such as the CC term
512 “nucleus” and the BP term “response to virus”. However, the only two predicted MF
513 terms have quite high consensus scores, “protein binding” and “protein
514 homodimerization activity”, suggesting that TRIM22 binds to itself to form a dimer.

515 We thus set out to test if human TRIM22 can form homodimer using
516 coimmunoprecipitation (**Figure 6A**). We first expressed Flag or GFP-tagged human
517 TRIM22 protein by co-transfecting the two plasmids into HEK293T cells. After 48-hour
518 incubation, cells were harvested and lysed (20 mM Tris-HCl pH7.5, 150 mM NaCl, 1
519 mM EDTA, 1% NP-40 with proteinase inhibitor). Both Flag and GFP-tagged TRIM22
520 were detected in cell lysate by western blot (Figure 6B). We then pulled down Flag-
521 tagged TRIM22 from the cell lysate. For each sample, 25 μ l protein A/G beads were
522 incubated with 1 μ g Flag antibody at 4 degree. Mouse IgG was used as a negative
523 control. After 2 hours, beads were washed with lysis buffer and then incubated with 500
524 μ g cell lysate at 4 degree for another 2 hours. Western blot showed that when Flag-tagged
525 TRIM22 was pulled down, GFP-tagged TRIM22 can be detected by GFP antibody
526 (Figure 6C), which showed that GFP-tagged and Flag-tagged TRIM22 bind together in
527 HEK293T. To further confirm this binding, we did co-IP in the opposite way. Flag-
528 tagged TRIM22 was also detected in immunoprecipitation of GFP-tagged TRIM22
529 (Figure 6D). These results reveal that TRIM22 can bind to itself, which is most likely to
530 form the homodimer.

531

532 **Conclusion**

533 In this work, we developed QAUST, a method to predict biological functions of protein
534 molecules using three main features: global and local protein structure similarity, protein-
535 protein interaction and functional sequence motifs. In our method, we constructed the 3D
536 structure from the amino acid sequence using I-TASSER. Functional analogs are then
537 identified by performing global and local structural similarity search through the

538 functional libraries, with the scoring function involving the confidence score of structural
539 predictions, sequence and structural similarity of the I-TASSER model with the
540 functional templates, and the local active site matches. We have also tried to improve the
541 performance of GO prediction by incorporating protein-protein interaction information,
542 especially in order to improve the prediction of GO terms under BP and CC aspects. We
543 further developed a novel predictor that extracts functional motifs that are related to a
544 specific GO term and used it as our third predictor.

545 On a set of 500 non-redundant proteins, QAUST is shown to have higher function
546 prediction accuracy than all other competing methods on most prediction tasks. This
547 performance advantage is mainly a result of combining three different predictors which
548 cover major aspects of proteins. Additionally, our three prediction components
549 complement each other in the sense that they contribute differently to the prediction of
550 the three aspects of GO. While PPI information improves significantly the prediction of
551 BP and CC terms, functional motifs detection is mainly useful in improving MF term
552 prediction. However, QAUST has a number of limitations that give room for possible
553 improvement in the future. One main limitation is that QAUST is much more expensive
554 in terms of running time compared to the other methods as reported in Section S10 and
555 Table S4. The second limitation is that our method cannot be directly used to infer
556 functions that are not included in EC or GO systems since it solely infers protein
557 functions from existing protein annotations. Finally, given that the three components we
558 used work differently in predicting different aspects of GO, it may be helpful to weight
559 their scores differently depending on the nature of the GO term evaluated instead of
560 combining the scores in a simple consensus. In particular, advanced machine learning
561 methods such as deep learning [57–61], could help weight and combine the scores in a
562 more efficient way to obtain better prediction results which could be a possible future
563 improvement of this work.

564

565 **Authors' contributions**

566 FS designed and carried out the experiments and drafted the manuscript. AR drafted the
567 initial manuscript and provided the EC data set and template libraries. MA helped design
568 the webserver. STA revised the manuscript. SM and PSH helped draft the initial

569 manuscript. ST and WC designed and conducted experimental validation. XG supervised
570 the algorithm design and the experiments, and helped draft and revise the manuscript.

571

572 **Competing interests**

573 The authors have declared no competing interest.

574

575 **Acknowledgement**

576 We thank Mr. Chengxin Zhang, Dr. Wei Zhang and Professor Yang Zhang for helpful
577 discussions. The research reported in this publication was supported by the King
578 Abdullah University of Science and Technology (KAUST) Office of Sponsored Research
579 (OSR) under Award No. URF/1/1976-04 and URF/1/1976-06. This work used the
580 Extreme Science and Engineering Discovery Environment (XSEDE) [62], which is
581 supported by National Science Foundation grant number ACI-1053575.

582

583 **References**

584 [1] U. Consortium, "UniProt: a hub for protein information," *Nucleic acids research*,
585 vol. 43, no. D1, pp. D204-D212, 2014.

586 [2] E. Boutet *et al.*, "UniProtKB/Swiss-Prot, the Manually Annotated Section of the
587 UniProt KnowledgeBase: How to Use the Entry View," (in English), *Plant*
588 *Bioinformatics: Methods and Protocols, 2nd Edition*, vol. 1374, pp. 23-54, 2016.

589 [3] S. F. Altschul *et al.*, "Gapped BLAST and PSI-BLAST: a new generation of
590 protein database search programs," (in English), *Nucleic Acids Research*, vol. 25, no. 17,
591 pp. 3389-3402, Sep 1 1997.

592 [4] S. R. Eddy, "Profile hidden Markov models," (in English), *Bioinformatics*, vol.
593 14, no. 9, pp. 755-763, 1998.

594 [5] J. Soding, "Protein homology detection by HMM-HMM comparison," (in
595 English), *Bioinformatics*, vol. 21, no. 7, pp. 951-960, Apr 1 2005.

- 596 [6] K. Karplus, C. Barrett, and R. Hughey, "Hidden Markov models for detecting
597 remote protein homologies," (in English), *Bioinformatics*, vol. 14, no. 10, pp. 846-856,
598 1998.
- 599 [7] W. Tian and J. Skolnick, "How well is enzyme function conserved as a function
600 of pairwise sequence identity?," *Journal of molecular biology*, vol. 333, no. 4, pp. 863-
601 882, 2003.
- 602 [8] B. Rost, "Enzyme function less conserved than anticipated," (in English), *Journal*
603 *of Molecular Biology*, vol. 318, no. 2, pp. 595-608, Apr 26 2002.
- 604 [9] R. D. Finn *et al.*, "InterPro in 2017-beyond protein family and domain
605 annotations," (in English), *Nucleic Acids Research*, vol. 45, no. D1, pp. D190-D199, Jan
606 4 2017.
- 607 [10] D. A. D. L. Morais *et al.*, "SUPERFAMILY 1.75 including a domain-centric gene
608 ontology method," (in English), *Nucleic Acids Research*, vol. 39, pp. D427-D434, Jan
609 2011.
- 610 [11] R. Rentzsch and C. A. Orengo, "Protein function prediction using domain
611 families," (in English), *Bmc Bioinformatics*, vol. 14, Feb 28 2013.
- 612 [12] G. Lopez, A. Rojas, M. Tress, and A. Valencia, "Assessment of predictions
613 submitted for the CASP7 function prediction category," (in English), *Proteins-Structure*
614 *Function and Bioinformatics*, vol. 69, pp. 165-174, 2007.
- 615 [13] Y. Zhang, "Protein structure prediction: when is it useful?," (in English), *Current*
616 *Opinion in Structural Biology*, vol. 19, no. 2, pp. 145-155, Apr 2009.
- 617 [14] J. Skolnick, J. S. Fetrow, and A. Kolinski, "Structural genomics and its
618 importance for gene function analysis," (in English), *Nature Biotechnology*, vol. 18, no.
619 3, pp. 283-287, Mar 2000.
- 620 [15] P. Aloy, E. Querol, F. X. Aviles, and M. J. E. Sternberg, "Automated structure-
621 based prediction of functional sites in proteins: Applications to assessing the validity of
622 inheriting protein function from homology in genome annotation and to protein docking,"
623 (in English), *Journal of Molecular Biology*, vol. 311, no. 2, pp. 395-408, Aug 10 2001.
- 624 [16] A. Roy, N. Srinivasan, and V. S. Gowri, "Molecular and structural basis of drift in
625 the functions of closely-related homologous enzyme domains: implications for function

626 annotation based on homology searches and structural genomics," *In silico biology*, vol.
627 9, no. 1, 2, pp. S41-S55, 2009.

628 [17] P. Bork, C. Sander, and A. Valencia, "Convergent Evolution of Similar
629 Enzymatic Function on Different Protein Folds - the Hexokinase, Ribokinase, and
630 Galactokinase Families of Sugar Kinases," (in English), *Protein Science*, vol. 2, no. 1, pp.
631 31-40, Jan 1993.

632 [18] R. V. Spriggs, P. J. Artymiuk, and P. Willett, "Searching for patterns of amino
633 acids in 3D protein structures," (in English), *Journal of Chemical Information and
634 Computer Sciences*, vol. 43, no. 2, pp. 412-421, Mar-Apr 2003.

635 [19] K. Kinoshita and H. Nakamura, "Identification of protein biochemical functions
636 by similarity search using the molecular surface database eF-site," (in English), *Protein
637 Science*, vol. 12, no. 8, pp. 1589-1595, Aug 2003.

638 [20] D. T. H. Chang, C. Y. Chen, W. C. Chung, Y. J. Oyang, H. F. Juan, and H. C.
639 Huang, "ProteMiner-SSM: a web server for efficient analysis of similar protein tertiary
640 substructures," (in English), *Nucleic Acids Research*, vol. 32, pp. W76-W82, Jul 1 2004.

641 [21] P. F. Gherardini and M. Helmer-Citterich, "Structure-based function prediction:
642 approaches and applications," *Briefings in functional genomics & proteomics*, vol. 7, no.
643 4, pp. 291-302, 2008.

644 [22] R. A. Laskowski, J. D. Watson, and J. M. Thornton, "Protein function prediction
645 using local 3D templates," *Journal of molecular biology*, vol. 351, no. 3, pp. 614-626,
646 2005.

647 [23] C. Zhang, W. Zheng, P. L. Freddolino, and Y. Zhang, "MetaGO: Predicting Gene
648 Ontology of Non-homologous Proteins Through Low-Resolution Protein Structure
649 Prediction and Protein-Protein Network Mapping," *Journal of molecular biology*, 2018.

650 [24] A. Roy, J. Y. Yang, and Y. Zhang, "COFACTOR: an accurate comparative
651 algorithm for structure-based protein function annotation," (in English), *Nucleic Acids
652 Research*, vol. 40, no. W1, pp. W471-W477, Jul 2012.

653 [25] C. Brun, F. Chevenet, D. Martin, J. Wojcik, A. Guenoche, and B. Jacq,
654 "Functional classification of proteins for the prediction of cellular function from a
655 protein-protein interaction network," (in English), *Genome Biology*, vol. 5, no. 1, 2004.

- 656 [26] H. N. Chua, W. K. Sung, and L. Wong, "Using indirect protein interactions for
657 the prediction of Gene Ontology functions," (in English), *Bmc Bioinformatics*, vol. 8,
658 2007.
- 659 [27] A. L. Barabasi, N. Gulbahce, and J. Loscalzo, "Network medicine: a network-
660 based approach to human disease," (in English), *Nature Reviews Genetics*, vol. 12, no. 1,
661 pp. 56-68, Jan 2011.
- 662 [28] R. Magnez, B. Thiroux, S. Taront, Z. Segaula, B. Quesnel, and X. Thuru, "PD-
663 1/PD-L1 binding studies using microscale thermophoresis," *Scientific reports*, vol. 7, no.
664 1, p. 17623, 2017.
- 665 [29] L. Lan, N. Djuric, Y. Guo, and S. Vucetic, "MS-k NN: protein function prediction
666 by integrating multiple data sources," in *BMC bioinformatics*, 2013, vol. 14, no. 3, p. S8:
667 BioMed Central.
- 668 [30] R. You, Z. Zhang, Y. Xiong, F. Sun, H. Mamitsuka, and S. Zhu, "GOLabeler:
669 improving sequence-based large-scale protein function prediction by learning to rank,"
670 *Bioinformatics*, vol. 1, p. 9, 2018.
- 671 [31] M. Kulmanov, M. A. Khan, and R. Hoehndorf, "DeepGO: predicting protein
672 functions from sequence and interactions using a deep ontology-aware classifier,"
673 *Bioinformatics*, vol. 34, no. 4, pp. 660-668, 2017.
- 674 [32] Q. Gong, W. Ning, and W. Tian, "GoFDR: a sequence alignment based method
675 for predicting protein functions," *Methods*, vol. 93, pp. 3-14, 2016.
- 676 [33] Z. Zou, S. Tian, X. Gao, and Y. Li, "mldeepr: Multi-functional enzyme function
677 prediction with hierarchical multi-label deep learning," *Frontiers in genetics*, vol. 9,
678 2018.
- 679 [34] X. Gao, D. Bu, J. Xu, and M. Li, "Improving consensus contact prediction via
680 server correlation reduction," *BMC structural biology*, vol. 9, no. 1, p. 28, 2009.
- 681 [35] P. Chen *et al.*, "A sequence-based dynamic ensemble learning system for protein
682 ligand-binding site prediction," *IEEE/ACM transactions on computational biology and*
683 *bioinformatics*, vol. 13, no. 5, pp. 901-912, 2016.
- 684 [36] P. Chen, J. Z. Huang, and X. Gao, "LigandRFs: random forest ensemble to
685 identify ligand-binding residues from sequence information alone," in *BMC*
686 *bioinformatics*, 2014, vol. 15, no. 15, p. S4: BioMed Central.

687 [37] P. Chen, J. Li, L. Wong, H. Kuwahara, J. Z. Huang, and X. Gao, "Accurate
688 prediction of hot spot residues through physicochemical characteristics of amino acid
689 sequences," *Proteins: Structure, Function, and Bioinformatics*, vol. 81, no. 8, pp. 1351-
690 1362, 2013.

691 [38] Y. Zhang, "I-TASSER server for protein 3D structure prediction," *BMC*
692 *bioinformatics*, vol. 9, no. 1, p. 1, 2008.

693 [39] C. Zhang, P. L. Freddolino, and Y. Zhang, "COFACTOR: improved protein
694 function prediction by combining structure, sequence and protein-protein interaction
695 information," *Nucleic acids research*, p. gkx366, 2017.

696 [40] D. Szklarczyk *et al.*, "STRING v10: protein-protein interaction networks,
697 integrated over the tree of life," (in English), *Nucleic Acids Research*, vol. 43, no. D1, pp.
698 D447-D452, Jan 28 2015.

699 [41] E. C. Webb, *Enzyme nomenclature 1992. Recommendations of the Nomenclature*
700 *Committee of the International Union of Biochemistry and Molecular Biology on the*
701 *Nomenclature and Classification of Enzymes* (no. Ed. 6). Academic Press, 1992.

702 [42] M. Ashburner *et al.*, "Gene Ontology: tool for the unification of biology," (in
703 English), *Nature Genetics*, vol. 25, no. 1, pp. 25-29, May 2000.

704 [43] A. Roy, J. Yang, and Y. Zhang, "COFACTOR: an accurate comparative
705 algorithm for structure-based protein function annotation," *Nucleic acids research*, vol.
706 40, no. W1, pp. W471-W477, 2012.

707 [44] S. T. Wu and Y. Zhang, "LOMETS: A local meta-threading-server for protein
708 structure prediction," (in English), *Nucleic Acids Research*, vol. 35, no. 10, pp. 3375-
709 3382, May 2007.

710 [45] Y. X. Jiang *et al.*, "An expanded evaluation of protein function prediction
711 methods shows an improvement in accuracy," (in English), *Genome Biology*, vol. 17, Sep
712 7 2016.

713 [46] D. Piovesan, M. Giollo, E. Leonardi, C. Ferrari, and S. C. E. Tosatto, "INGA:
714 protein function prediction combining interaction networks, domain assignments and
715 sequence similarity," (in English), *Nucleic Acids Research*, vol. 43, no. W1, pp. W134-
716 W140, Jul 1 2015.

717 [47] C. T. Porter, G. J. Bartlett, and J. M. Thornton, "The Catalytic Site Atlas: a
718 resource of catalytic sites and residues identified in enzymes using structural data," (in
719 English), *Nucleic Acids Research*, vol. 32, pp. D129-D133, Jan 1 2004.

720 [48] Y. Zhang and J. Skolnick, "TM-align: a protein structure alignment algorithm
721 based on the TM-score," (in English), *Nucleic Acids Research*, vol. 33, no. 7, pp. 2302-
722 2309, 2005.

723 [49] U. Consortium, "The universal protein resource (UniProt) in 2010," *Nucleic acids
724 research*, vol. 38, no. suppl_1, pp. D142-D148, 2009.

725 [50] Y. Li *et al.*, "DEEPre: sequence-based enzyme EC number prediction by deep
726 learning," *Bioinformatics*, vol. 34, no. 5, pp. 760-769, 2017.

727 [51] P. Radivojac *et al.*, "A large-scale evaluation of computational protein function
728 prediction," *Nature methods*, vol. 10, no. 3, p. 221, 2013.

729 [52] T. C. Mueser, C. E. Jones, N. G. Nossal, and C. C. Hyde, "Bacteriophage T4 gene
730 59 helicase assembly protein binds replication fork DNA. The 1.45 Å resolution crystal
731 structure reveals a novel α -helical two-domain fold1," *Journal of molecular biology*, vol.
732 296, no. 2, pp. 597-612, 2000.

733 [53] S. D. Barr, J. R. Smiley, and F. D. Bushman, "The interferon response inhibits
734 HIV particle production by induction of TRIM22," *PLoS pathogens*, vol. 4, no. 2, p.
735 e1000007, 2008.

736 [54] A. Di Pietro *et al.*, "TRIM22 inhibits influenza A virus infection by targeting the
737 viral nucleoprotein for degradation," *Journal of virology*, vol. 87, no. 8, pp. 4523-4533,
738 2013.

739 [55] C. Yang *et al.*, "Interferon alpha (IFN α)-induced TRIM22 interrupts HCV
740 replication by ubiquitinating NS5A," *Cellular & molecular immunology*, vol. 13, no. 1, p.
741 94, 2016.

742 [56] J. Lou, Y. Wang, X. Zheng, and W. Qiu, "TRIM22 regulates macrophage
743 autophagy and enhances Mycobacterium tuberculosis clearance by targeting the nuclear
744 factor–multiplicity κ B/beclin 1 pathway," *Journal of cellular biochemistry*, vol. 119, no.
745 11, pp. 8971-8980, 2018.

746 [57] Z. Xia *et al.*, "DeeReCT-PolyA: a robust and generic deep learning method for
747 PAS identification," 2018.

748 [58] R. Umarov, H. Kuwahara, Y. Li, X. Gao, and V. Solovyev, "Promoter analysis
749 and prediction in the human genome using sequence-based deep learning models,"
750 *Bioinformatics*, vol. 1, p. 8, 2019.

751 [59] J.-S. Kim, X. Gao, and A. Rzhetsky, "RIDDLE: Race and ethnicity Imputation
752 from Disease history with Deep LEarning," *PLoS computational biology*, vol. 14, no. 4,
753 p. e1006106, 2018.

754 [60] Y. Li *et al.*, "Dlbi: deep learning guided bayesian inference for structure
755 reconstruction of super-resolution fluorescence microscopy," *Bioinformatics*, vol. 34, no.
756 13, pp. i284-i294, 2018.

757 [61] Y. Li, C. Huang, L. Ding, Z. Li, Y. Pan, and X. Gao, "Deep learning in
758 bioinformatics: introduction, application, and perspective in big data era," *arXiv preprint*
759 *arXiv:1903.00342*, 2019.

760 [62] J. Towns *et al.*, "XSEDE: Accelerating Scientific Discovery," (in English),
761 *Computing in Science & Engineering*, vol. 16, no. 5, pp. 62-74, Sep-Oct 2014.

762

763

764 **Figure legends**

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

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

769

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

771 The query sequence is searched against UniRef90 database by PSI-BLAST to identify
772 sequence homologs with GO term annotation. For a GO term of interest, λ , the identified
773 homologs are divided into two sets: the “annotated set” (purple) which contains
774 homologs annotated with λ , and “not-annotated set” (green) which consists of homologs
775 not associated with λ . From each of the two sets, frequent motifs, i.e. continuous
776 sequence fragments, are extracted. For illustration purposes, only three five-residue-long
777 motifs from each set are drawn. The GO term λ is predicted with confidence score
778 $n_q(\lambda)/N(\lambda) \cdot [1 - n_q(\lambda^c)/N(\lambda^c)]$. Here, $N(\lambda)$ and $N(\lambda^c)$ are the total number of
779 extracted frequent motifs for “annotated set” and “not-annotated set”, correspondingly;
780 while $n_q(\lambda)$ and $n_q(\lambda^c)$ are the number of frequent motifs from “annotated set” and “not-
781 annotated set” that match the query sequence, respectively. In this example, only the
782 motif “CLPFD” from “annotated set” matches the query, making the confidence score
783 equals to $1/3 \cdot [1 - 0/3] = 1/3$.

784

785 **Figure 3 Precision-recall curves for EC prediction by QAUST, COFACTOR,**
786 **LOMETS, HHsearch and BLAST**

787

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

789 GO prediction performance of our method based on different sets of features, and seven
790 other methods for each of the three GO branches.

791

792 **Figure 5 A study case for protein function prediction using QAUST**

793 **A.** The superimposition between the query (P13342, in cyan) and the closest template in
794 the database (PDB ID 2CT8A, in magenta) based on the structural alignment generated

795 by TM-align. **B.** Predicted BP terms for protein P13342. The six BP terms (the root term,
796 Biological Process, is a naïve term, which is not counted) shown are the experimentally
797 annotated terms. The colored contours represent the BP terms that are predicted by the
798 corresponding methods.

799

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

801 **A.** Illustration of the coimmunoprecipitation method to validate the homodimerization of
802 TRIM22. We expressed Flag- or GFP-tagged human TRIM22 protein by co-transfecting
803 two plasmid into HEK293T cells. If TRIM22 forms a homodimer, when Flag-tagged
804 TRIM22 or GFP-tagged TRIM22 is pulled down, both Flag-tagged and GFP-tagged
805 TRIM22 should be detected by the corresponding antibodies (4 combinations in total). **B.**
806 Both Flag- and GFP-tagged TRIM22 expressed in HEK293T cells, detected by Western
807 Blot. **C.** For Flag-immunoprecipitation, both Flag-tagged and GFP-tagged TRIM22 are
808 detected by the corresponding antibodies, whereas mouse IgG is used as a negative
809 control. **D.** For GFP-immunoprecipitation, both Flag-tagged and GFP-tagged TRIM22
810 are detected by the corresponding antibodies, whereas rabbit IgG is the negative control.

811

812 **Tables**

813 **Table 1 Accuracy values of EC prediction for QAUST as well as five other methods**

814 **Table 2 Fmax values of each branch of GO for QAUST as well as prediction**
815 **methods**

816 **Table 3 Smin (minimum semantic distance) values of each branch of GO for**
817 **QAUST as well as other prediction methods**

818

819 **Supplementary materials**

820 **Figure S1 Precision-recall curves for GO prediction using different structure**
821 **methods**

822 Precision-recall curves for GO prediction by QAUST's structure similarity based pipeline
823 (I-TASSER + $P_{\text{structure}}$), alternative implementations using low resolution homology

824 models (HHsearch + $P_{\text{structure}}$ and LOMETS + $P_{\text{structure}}$) and baseline algorithms (HHsearch
825 and LOMETS).

826

827 **Figure S2 GO prediction performance of QAUST compared to the prediction**
828 **performance when using PPI and motif features only for each one of the three GO**
829 **branches**

830

831 **Figure S3 Set of predicted MF GO terms for protein P01574**

832

833 **Figure S4 Set of predicted CC GO terms for protein P32157**

834

835 **Table S1 P values from the Mann-Whitney U test to assess improvement/decrease**
836 **of QAUST performance in GO prediction compared to other methods**

837

838 **Table S2 F_{max} values from GO prediction using only experimental GO annotations**
839 **for motif detection (QAUST compared to using both experimental and predicted**
840 **terms)**

841

842 **Table S3 Predicted GO functions for protein TRIM22 with the obtained confidence**
843 **score for each one of the three GO branches**

844

845 **Table S4 Running time of QAUST and other prediction methods for GO prediction**
846 **on our dataset (500 sequences)**