

1 **Towards a systematic characterization of protein complex function: a** 2 **natural language processing and machine-learning framework**

3
4 Varun S. Sharma^{1,2}, Andrea Fossati^{3,4}, Rodolfo Ciuffa¹, Marija Buljan⁵, Evan G. Williams⁶, Zhen Chen⁷,
5 Wenguang Shao¹, Patrick G.A. Pedrioli¹, Anthony W. Purcell⁸, María Rodríguez Martínez⁹, Jiangning
6 Song^{8,10,*}, Matteo Manica^{9,*}, Ruedi Aebersold^{1, 11,*}, and Chen Li^{1,8,12,*}

7
8 ¹Department of Biology, Institute of Molecular Systems Biology, ETH Zürich, Switzerland
9 ²Institute for Neurodegenerative Diseases, University of California, San Francisco, CA 94143, USA
10 ³Quantitative Biosciences Institute (QBI) and Department of Cellular and Molecular Pharmacology,
11 University of California, San Francisco, CA 94158, USA
12 ⁴J. David Gladstone Institutes, San Francisco, CA 94158, USA
13 ⁵Empa - Swiss Federal Laboratories for Materials Science and Technology, St. Gallen, Switzerland
14 ⁶Luxembourg Centre for Systems Biomedicine, University of Luxembourg, Esch-sur-Alzette,
15 Luxembourg
16 ⁷Collaborative Innovation Center of Henan Grain Crops, Henan Agricultural University, Zhengzhou
17 450046, China
18 ⁸Department of Biochemistry and Molecular Biology and Monash Biomedicine Discovery Institute,
19 Monash University, Clayton, Victoria, Australia
20 ⁹IBM Research Europe, Zürich, Switzerland
21 ¹⁰Monash Centre for Data Science, Faculty of Information Technology, Monash University, Melbourne,
22 VIC 3800, Australia
23 ¹¹Faculty of Science, University of Zürich, Switzerland.
24 ¹²Lead contact
25
26 *Correspondence: Jiangning.Song@monash.edu (J. Song), TTE@zurich.ibm.com (M. Manica),
27 aegersold@imsb.biol.ethz.ch (R. Aebersold), and Chen.Li@monash.edu (C. Li).

28 **Summary**

29 It is a general assumption of molecular biology that the ensemble of expressed molecules, their activities
30 and interactions determine biological processes, cellular states and phenotypes. Quantitative abundance
31 of transcripts, proteins and metabolites are now routinely measured with considerable depth via an array
32 of “OMICS” technologies, and recently a number of methods have also been introduced for the parallel
33 analysis of the abundance, subunit composition and cell state specific changes of protein complexes. In
34 comparison to the measurement of the molecular entities in a cell, the determination of their function
35 remains experimentally challenging and labor-intensive. This holds particularly true for determining the
36 function of protein complexes, which constitute the core functional assemblies of the cell. Therefore, the
37 tremendous progress in multi-layer molecular profiling has been slow to translate into increased
38 functional understanding of biological processes, cellular states and phenotypes. In this study we
39 describe PCfun, a computational framework for the systematic annotation of protein complex function
40 using Gene Ontology (GO) terms. This work is built upon the use of word embedding— natural language
41 text embedded into continuous vector space that preserves semantic relationships— generated from the
42 machine reading of 1 million open access PubMed Central articles. PCfun leverages the embedding for
43 rapid annotation of protein complex function by integrating two approaches: (1) an unsupervised
44 approach that obtains the nearest neighbor (NN) GO term word vectors for a protein complex query
45 vector, and (2) a supervised approach using Random Forest (RF) models trained specifically for
46 recovering the GO terms of protein complex queries described in the CORUM protein complex database.
47 PCfun consolidates both approaches by performing the statistical test for the enrichment of the top NN
48 GO terms within the child terms of the predicted GO terms by RF models. Thus, PCfun amalgamates
49 information learned from the gold-standard protein-complex database, CORUM, with the unbiased
50 predictions obtained directly from the word embedding, thereby enabling PCfun to identify the potential
51 functions of putative protein complexes. The documentation and examples of the PCfun package are
52 available at <https://github.com/sharmavaruns/PCfun>. We anticipate that PCfun will serve as a useful tool
53 and novel paradigm for the large-scale characterization of protein complex function.

54 INTRODUCTION

55 Proteins are known to catalyze and control the vast majority of the reactions of cellular biochemistry
56 (Aebersold and Mann, 2016). Frequently they exert their function only if they stably interact in precise
57 stoichiometric ratios with other proteins in the form of complex macromolecular structures, a notion that
58 has been encapsulated in the term “modular cell biology” (Hartwell et al., 1999). With the advent of high
59 throughput ‘OMICS’ technologies for the study of complex biological systems, it is now possible to
60 accurately quantify and identify different types of biologically relevant molecules across various
61 conditions. However, determining the biological functions and phenotypes of these assemblies has
62 remained challenging and requires a functional understanding of the molecular functions of its
63 components and associations. Detailed biochemical and cell biological studies have identified the
64 composition and even the atomic structures of numerous protein complexes with well-defined roles in a
65 variety of fundamental biological processes (Hewick et al., 2003), such as in their participation in
66 transcriptional regulation (Aranda et al., 2015; Simonis et al., 2004; Tan et al., 2007; Webb and Westhead,
67 2009), cell cycle control (Becher et al., 2018; Chen et al., 2019; D'Avino et al., 2009) and signal
68 transduction (Pawson and Nash, 2000; Rebois and Hebert, 2003). Protein complexes can, therefore, be
69 considered essential agents and indicators of cellular functionality.

70 Recent technical advances, particularly in mass spectrometry (MS) - based proteomics have
71 greatly enhanced our capacity to determine the composition, stoichiometry and abundance of known
72 protein complexes and to identify new entities. These methodologies also support the systematic
73 identification of compositional or quantitative changes in complexes as a function of cellular state. These
74 include methods such as BF-MS (Biochemical Fractionation Mass Spectrometry) (Carlson et al., 2019;
75 Heusel et al., 2019; Heusel et al., 2020; Rosenberger et al., 2020; Stacey et al., 2017; Szklarczyk et al.,
76 2019), Affinity Purification MS (AP-MS), Cross-Linking MS (XL-MS) (Leitner et al., 2016; Leitner et
77 al., 2012; Liu et al., 2015) and limited proteolysis (LiP) (Schopper et al., 2017) and thermal proteome
78 profiling (TTP) (Mateus et al., 2020). Compared to the experimental detection of new protein complexes,
79 the determination of their biochemical or cellular function has significantly lagged behind because
80 experimentation with specific complexes is highly challenging. Given the challenge of characterizing

81 the function of protein complexes, hypotheses regarding the functional roles in which a newly discovered
82 protein complex participates are typically generated by careful manual review of prior literature.

83 The standard approach to manual literature review for identifying the putative function of a
84 protein complex consists of the search for publications and database entries about the individual protein
85 subunits and followed by the consolidation of the retrieved information. However, this manual curation
86 presents several limitations that can make it highly inefficient and highly biased. First, exhaustive
87 literature curation for all proteins belonging to even a single complex can easily become prohibitively
88 time-consuming due to the sheer volume of publications required to parse through. Given that the manual
89 curation for retrieving high-confidence functional annotations of a single protein complex can be
90 extremely laborious, performing such annotation on dozens or hundreds of novel entities discovered in
91 large scale complex centric proteomic fractionation experiments quickly becomes prohibitive. Second,
92 although protein complex databases, such as CORUM (Giurgiu et al., 2019) and Complex Portal (Meldal
93 et al., 2019) offer experimentally and manually validated functional annotations for better-studied
94 protein complexes (i.e. the ground truth), the literature-based functional annotations for the same protein
95 complexes can be highly dissimilar across different databases. Third, some proteins are multifunctional
96 and may have unique roles in different protein complexes, thereby highlighting that the function of
97 protein complexes is not simply the aggregate of their subunits' functions (Jeffery, 2015; Matalon et al.,
98 2014; Nakabayashi et al., 2014). As a case in point, we conducted a preliminary examination of the Gene
99 Ontology (GO) terms annotated for whole protein complexes in the CORUM database compared to each
100 individual subunit's GO term annotations in the QuickGO database (Binns et al., 2009). The results
101 showed that 2155 (61.4%), 319 (9.1%), and 169 (4.8%) heteromeric protein complexes in CORUM
102 contained at least one novel biological process, molecular function, or cellular component term that was
103 not annotated for any individual subunit's QuickGO entry, respectively. In other words, certain proteins
104 may participate in emergent functionality when assembled in a macromolecular complex that would be
105 non-obvious based on the known functions of the individual protein complex's subunits. We therefore
106 argue that it is of great importance to employ computational techniques that can assist large-scale

107 predictions of protein complex functions and provide useful insights and guidance for the follow-up
108 functional characterization experiments.

109 Given that the nature of information encoded in natural language-based functional descriptions
110 of protein complexes is fundamentally unstructured, it is challenging for traditional bioinformatic and
111 data mining approaches to meaningfully distill the information from specific publications. However,
112 computational methods from text-mining — the field concerned with computationally extracting
113 information from unstructured natural language text— have been also successfully applied to a variety
114 of biomedical problems and provide a promising avenue to address our task. These include extraction of
115 protein-protein relations and functions (Islamaj Dogan et al., 2019; Li et al., 2019b; Manica et al., 2019;
116 Subramani et al., 2015; Yu et al., 2018), determination of protein structure (Gaizauskas et al., 2003),
117 protein localization (Cejuela et al., 2018), and gene-disease relationships (Pletscher-Frankild et al., 2015).
118 However, to the best of our knowledge, no computational tool designed specifically for annotating the
119 functions of protein complexes has been described to date.

120 To address this dearth in direct functional annotation methods for protein complexes, in this work
121 we integrate text-mining and machine-learning techniques into a hybrid computational framework,
122 termed PCfun, that can be applied to large scale complex-centric proteome experiments for predicting
123 the function of protein complexes. At a high level, PCfun is developed based upon word embedding
124 generated from the machine reading of >1 million open access PubMed Central (PMC) articles, whereby
125 both unsupervised and supervised machine learning algorithms were used to generate two separate lists
126 of predicted functional Gene Ontology (GO; biological process, molecular function and cellular
127 component) terms for a queried protein complex. Following, a supervised machine-learning model
128 trained on the associations between protein complexes and their GO terms documented in the CORUM
129 database was used to predict a second list of candidate functional terms. Hence, the unsupervised
130 candidate list provides functional predictions solely based upon the word vector relationships observed
131 within the embedding that are unbiased to protein complex-function associations while the supervised
132 candidate list tailors the annotations to relationships similar to the CORUM database. In order to leverage
133 the insights provided by both approaches we attempted to consolidate the two lists by leveraging the

134 hierarchical structure of the Gene Ontology by testing for enrichment of certain supervised terms within
135 the unsupervised list. An adapted leave-one-out cross-validation scheme was used to test the system's
136 performance and suggested that PCfun achieved outstanding prediction performance with AUC values
137 of 0.895, 0.927, and 0.957 for biological process, molecular function, and cellular component terms,
138 respectively. In addition, we compared the prediction outcomes by PCfun and the GO annotations from
139 the Complex Portal database (Meldal et al., 2019) using protein complexes not documented in the
140 CORUM database. For the biological function and cellular component categories, PCfun predicted
141 similar (i.e. semantic similarity ≥ 0.5) GO terms that covered more than half of the Complex Portal's
142 ground-truth annotations for 52.8% and 69.7% of the protein complexes in the biological process and
143 cellular component categories respectively. In contrast, the molecular function category only achieved
144 12.9% coverage, which might be explained by the lowest similarity between the CORUM and Complex
145 Portal annotated GO terms for this category. Taken together, we anticipate that PCfun will serve as an
146 accurate annotation tool for protein complex function and increase our better understanding of the
147 functional roles of protein complexes in biological systems.

148

149 **RESULTS**

150 **The Architecture of PCfun for Predicting the Function of Protein Complexes**

151 PCfun development consisted of two main steps. The first step, as shown in **Figure 1A**, was based on
152 the building of the word embedding itself. Approximately 1 million open access articles were
153 downloaded from the PubMed Central Repository and their full texts were processed (STAR Methods)
154 as described in Manica *et al.* (Manica et al., 2019) to populate a text corpus. After consolidation of the
155 text corpus, the fastText implementation of the skip-gram context prediction embedding algorithm was
156 employed onto the text corpus (STAR Methods) to construct a word-embedding representation of the
157 text: 500-dimensional, continuous real-valued vector representations based on the subwords extracted
158 from the corpus. These word vectors were constructed for character n -grams allowing for the creation of
159 a word vector for any natural language query (see **Figure 2A** for a graphic illustration of the word
160 vectors). Using this property of character n -gram embeddings, we next extracted sub-embeddings

161 consisting of all protein complex and GO term (split into biological process, molecular function, and
162 cellular component classes) queries. As a result of this first step (**Figure 1A**) we obtained five-word
163 vector sub-embeddings corresponding to protein complex queries (sub-embedding for each naming
164 scheme used for a protein complex name: canonical *vs.* subunit name) and GO term queries (three sub-
165 embeddings corresponding to biological process, molecular function, and cellular component terms).

166 The next step was then to construct the models for functional annotation of protein complex
167 queries. As depicted in **Figures 1B** and **C**, we employed two strategies capable of returning ranked
168 protein complex – GO associations: (i) an unsupervised nearest-neighbor approach illustrated in **Figure**
169 **1B**, and (ii) a supervised machine learning approach displayed in **Figure 1C**. The first algorithm is
170 agnostic to the question of functional annotation of protein complexes and was based solely upon
171 contextual relationships. In contrast, the second algorithm is a tailored approach trained specifically to
172 recover functional terms for a protein complex query. The rationale for using two distinct approaches is
173 that they are likely to produce complementary, and potentially, if combined, more informative outputs.

174 As seen in **Figure 1B**, we built a k -d tree (k -dimensional tree) (Freidman et al., 1977), a space-
175 partitioning structure for storing the sub-embeddings' vectors of GO terms to enable rapid application
176 of a nearest neighbor algorithm that was able to shortlist GO terms ranked by cosine similarity between
177 the queried protein complex vector and each word vector for a GO term, for recovering CORUM's
178 ground-truth annotations of each protein complex. The supervised machine-learning models (**Figure 1C**),
179 on the other hand, learned from the experimentally verified protein complex – GO term associations and
180 were therefore able to accurately cover these ground-truth in the CORUM database. We constructed and
181 evaluated four widely applied machine-learning algorithms: RF (Breiman, 2001), Logistic Regression
182 (LR) (Lecessie and Vanhouwelingen, 1992; Yu et al., 2011), and Naïve Bayes (NB) classifier with
183 Gaussian distribution and a Bernoulli distribution (Zhang, 2004) classifiers. A ranked list of GO term
184 annotations was generated by both the unsupervised k -d tree algorithm (**Figure 1B**) and the supervised
185 machine learning models (**Figure 1C**).

186 Finally, to combine the prediction outcomes from the RF and the k -d tree, a hypergeometric test
187 was conducted to test for the functional enrichment of k -d tree terms within the child nodes of each RF

188 predicted term (**Figure 1D**). A visualization of a GO direct acyclic graph (DAG) structure for
189 functionally enriched predicted GO terms was performed to represent the contextual information of
190 predicted GO terms of biological process, molecular function and cellular component, respectively
191 (**Figure 1E**). Given a protein complex of interest, PCfun first applies the two models to generate two
192 prediction lists using k -d tree and RF, and then visualizes the prediction outcome via the functional
193 enrichment analysis and the GO DAG structure (STAR Methods).

194

195 **Benchmarking the prediction performance of PCfun**

196 We systematically evaluated the prediction performance of PCfun. We first separately assessed the
197 predictive ability of the unsupervised k -d tree and the RF model, for annotating protein complex
198 functions. Further, we independently compared the prediction outputs from the enrichment analysis of
199 PCfun with the functional annotations documented in the Complex Portal database (Meldal et al., 2019).

200

201 *The word-embedding and k -d tree facilitated the ranking of potential GO terms for protein complexes*

202 A useful property of word-embeddings is that words with related semantic meanings have corresponding
203 word vectors that exist closer to each other in the word vector space - as measured by the cosine similarity
204 (i.e. same orientation) – compared to words that have very different meanings. Therefore, one can find
205 words with a similar meaning to an input query word by simply finding the nearest neighbors of the input
206 query word vector. To aid in rapid nearest-neighbor calculations for these large sub-embeddings, we
207 stored each sub-embedding into a k -d tree, allowing us to efficiently retrieve word vectors that were
208 similar to the input query vector. To evaluate the quality of the word-embeddings, we performed
209 principal component analysis (PCA) of the word vectors for each extracted sub-embedding of different
210 types, including biological process vectors, molecular function vectors, cellular component vectors, and
211 the protein complex vectors with the two naming schemes (STAR Methods). **Figure 2B** demonstrates
212 that the sub-embeddings' word vectors of each type are well clustered, indicating the reliable quality of
213 the word-embedding.

214 We measured the ability of each GO term class sub-embedding to recover the ground-truth
215 functional annotations for a protein complex from CORUM by recording the number of nearest
216 neighbors (ranked by their cosine similarity) required to recover 100% of the ground-truth functional
217 annotations for an input protein complex query. We hypothesized that the results might change
218 depending on the name used to represent a protein complex. Additionally, considering that de novo
219 detected protein complexes will not be characterized with an accepted name, we proposed the subunit
220 naming scheme for a protein complex that would still allow for the functional annotation of even newly
221 identified protein complexes by PCfun. Therefore, we tested the two protein complex naming schemes'
222 sub-embeddings (STAR Methods). **Figures 2C, D** indicate that the sub-embeddings required on average
223 13487, 5119, 2692, and 11044, 5214, 1894 nearest neighbors to recover the ground truth for biological
224 process, molecular function, and cellular component categories using the canonical names and subunits
225 names, respectively. It is evident that in order to recover CORUM's ground-truth annotations, a large
226 number of nearest neighbors are required. We therefore subsequently performed manual literature search
227 based on the top nearest neighbor GO terms for certain protein complexes and observed that the predicted
228 GO terms were actually still quite informative and were recovering known biological knowledge.

229 We chose a representative example protein complex, namely "SMAD2-SMAD4-FAST1-TGIF-
230 HDAC1 complex, TGF (beta) induced", for a manual literature review comparison to the k -d tree nearest-
231 neighbor results. We used the subunit naming scheme ("smad4 tgif1 smad2 hdac1 foxh1") for generation
232 of its corresponding word vector and then queried the vector into the biological function, molecular
233 function and cellular component sub-embedding k -d trees. While the k -d trees required 27400, 2182 and
234 990 nearest neighbor terms in the biological function, molecular function and cellular component trees,
235 respectively, to recover the six CORUM annotated GO terms (DNA topological change; negative
236 regulation of transcription, DNA-templated; DNA binding; transforming growth factor beta receptor
237 signaling pathway; chromosome organization; nucleus) for this protein complex, the top returned k -d
238 tree nearest neighbors (**Supplemental Table S1**) still provided relevant GO terms that had related
239 biological meanings. For example, the top 10 nearest neighbors for the biological function category are
240 terms all related to the TGF β or bone morphogenic protein response. According to Massague *et al*

241 (Massague et al., 2005), the SMAD proteins accumulate in the nucleus to execute transcriptional control
242 in response to TGF β signal transduction and may be co-activated or co-repressed by various DNA-
243 binding co-factors. We observed that “negative regulation of Smad protein signal transduction” was the
244 8th nearest neighbor term for the queried protein complex vector, which recovers the role of the co-
245 repressor activity of HDAC1 and TGIF that act to repress the transcriptional control of the activated and
246 nuclear localized SMAD2:SMAD4 subcomplex (Liberati et al., 1999; Wicks et al., 2000). The ranking
247 of GO terms by the k -d tree for this protein complex is listed in **Supplementary Files 1-3**. In summary,
248 the top-nearest neighbor results of the k -d tree do provide insights into the relevant biology, but
249 demonstrate poor ability to recover CORUM’s ground-truth annotations. This suggests the necessity of
250 building supervised-learning models to systematically and statistically improve the predictive outcomes.

251
252 *Supervised machine-learning models greatly improved the performance of ground-truth recovery of GO*
253 *terms in CORUM*

254 In order to improve the performance of ground-truth recovery of CORUM, we implemented supervised
255 machine-learning classifiers based on the word vectors for a ‘protein complex-GO association’ pair
256 (termed ‘PC-GO’). In our study, the annotated association of a PC-GO term was regarded as a positive
257 sample, whereas randomly sampled synthetic pairs of a protein complex and other GO terms that were
258 not associated in CORUM were regarded as negative ones. As the negative samples significantly
259 outnumbered the positive samples in the resulting datasets, we generated five different training datasets
260 with randomly selected negative samples and all positives for each protein complex to ensure an equal
261 distribution of positive and negative samples for training the classifier (STAR Methods). This process
262 was conducted for both naming schemes. With the training datasets, we assessed the performance of
263 three machine-learning classification algorithms, including RF, LR, and NB with a Gaussian or Bernoulli
264 prior (NB_Gauss or NB_Bernoulli, respectively), through the adapted ‘protein complex’-leave-one-out
265 cross-validation strategy using standard performance measures (STAR Methods).

266 Across these classifiers, RF consistently performed the best as measured by all performance
267 metrics (**Figure 3, Supplementary Table S2, and Supplementary Figures S1 and S2**) and achieved a

268 robust performance across the two naming schemes. For example, via the ‘protein complex’-leave-one-
269 out cross-validation strategy, the RF classifier achieved AUC values of 0.885 and 0.895 for biological
270 function, 0.925 and 0.927 for molecular function, and 0.951 and 0.957 for cellular component category
271 for protein complexes with conventional and gene combinational names, respectively. In addition, we
272 also observed that the resulting GO term lists predicted by the RF classifiers were able to significantly
273 reduce the number of nearest neighbors needed to recover the majority of the ground-truth GO term
274 annotations for a protein complex when compared to the nearest-neighbor results from querying the k -d
275 tree (**Figure 3D**). For example, for protein complexes with subunit names, the RF classifier predicted
276 terms were able to recover 80.5%, 83.6%, and 89.2% of CORUM’s ground-truth in 102, 49, and 11
277 positively predicted terms for biological process, molecular function, and cellular component,
278 respectively.

279 While the RF classifiers performed well to recover the ground-truth as documented in the
280 CORUM database, performance of the supervised approach may belie the inherent bias to the database
281 that it was trained upon. Although protein complexes within CORUM have been extensively studied and
282 the GO term annotations have been manually curated, there is an extant right skew in the frequency of
283 GO terms with low to middle depth, based on the GO DAG structure. As shown in **Supplementary**
284 **Figure S3**, we observed that the logged frequency of a particular GO term (i.e. the number of times a
285 GO term has been annotated in CORUM) versus each GO term’s depth in the GO DAG structure reveals
286 a biased annotation distribution for GO terms in CORUM. For example, the biological process term
287 ‘Regulation of transcription DNA templated’, molecular function term ‘DNA binding’ and cellular
288 component term ‘Nucleus’ were annotated in 233, 278, and 702 protein complexes respectively out of
289 3511 total protein complexes in the CORUM database. Such over-annotated GO terms could bias
290 machine-learning algorithms in favor of selecting these highly abundant annotations. Therefore, to
291 address the biases of the dataset that the RF classifier was trained upon, we supplemented the predicted
292 terms from the RF classifier with the predicted nearest neighbors from the k -d tree. A graphical
293 illustration of the combination of the RF and k -d tree prediction lists together as well as an example is
294 shown in **Figure 4** (STAR Methods).

295

296 *Independent test demonstrates divergent GO term predictions by PCfun compared to Complex Portal*

297 We accessed the prediction performance of PCfun (trained on CORUM) with the testing data from a
298 different database, Complex Portal. We first identified that only 110 annotated human protein complexes
299 (i.e. complexes with identical subunits) were shared by CORUM and Complex Portal. We subsequently
300 interrogated the semantic similarity of the biological process, molecular function and cellular component
301 terms for these complexes between the two databases (STAR Methods). The heatmaps of the pairwise
302 similarity scores using the method reported by Wang *et al* (Wang et al., 2007) are shown in the left panel
303 of **Figure 5**. The average semantic similarity scores of biological process, molecular function and
304 cellular component categories were 0.40, 0.34 and 0.54, respectively, suggesting that even for complexes
305 with the same subunit composition the annotations of CORUM and Complex Portal are dissimilar. More
306 stringently, we examined the numbers of identical GO terms used for each overlapping protein complex.
307 As a result, less than one GO term across all categories (0.22 biological process term, 0.05 molecular
308 function term and 0.13 cellular component term, respectively) was shared on average per protein
309 complex between CORUM and Complex Portal. This means that the annotations for protein complexes
310 are extremely divergent across the two databases, making it challenging for PCfun (built on CORUM)
311 to accurately cover the GO annotations in Complex Portal. We then sought to gauge the approximate
312 similarity of GO terms predicted by PCfun with Complex Portal annotations by assessing the pairwise
313 semantic similarity (STAR Methods). The right panels of **Figure 5** show the comparison between
314 predicted biological process, molecular function and cellular component terms by PCfun and the
315 Complex Portal annotations for the non-overlapping protein complexes (i.e. with <50% of overlapping
316 subunits). For biological process as shown in **Figure 5A**, 15 (approximately 44.1%) non-overlapping
317 complexes covered 90-100% of similar terms (semantic similarity ≥ 0.5) compared to Complex Portal
318 biological process annotations, while 14 complexes (41.2%) had divergent predictions (i.e. coverage
319 between 0-10%) compared to the annotations in the Complex Portal. Similarly, for cellular component
320 (**Figure 5C**), approximately 69.7% (23) of the complexes covered 90-100% of similar terms compared
321 to Complex Portal cellular component annotations 30.3% (10). In contrast, PCfun demonstrated highly

322 divergent predictions for the molecular function category with only 1 (3.2%) complex covering 90-100%
323 of similar terms and 27 (87.1%) complexes demonstrating 0-10% coverage of similar terms when
324 compared to Complex Portal's annotations. The low coverage of PCfun prediction for the molecular
325 function category might be related to the limited overlap between the molecular function annotations
326 from CORUM and Complex Portal (the left panel of **Figure 5B**). It is noteworthy that different studies
327 and databases may have divergent annotations for a protein or protein complex. Despite the low semantic
328 similarities of all the biological process, molecular function and cellular component terms between the
329 two databases, PCfun still demonstrates its ability to accurately recover and reliably predict functions
330 accurately in the context of the database it was trained.

331

332 **DISCUSSION**

333 Advances in high throughput omics techniques have allowed for fast and deep identification and
334 quantification of functional macro-molecules in the cell. Understanding the functions of these molecules,
335 particularly for protein complexes, is a crucial step to understand and model biological activities in the
336 cell. However, annotating the function of a protein complex is a great challenge, given the huge and
337 occasionally contradictory volume of literature on each of its subunits. Currently, manual literature
338 review is still the most common way to annotate the functions of protein complexes in major databases,
339 such as CORUM and Complex Portal. Computational methods for gene function prediction (Barutcuoglu
340 et al., 2006; Guan et al., 2008; Stojanova et al., 2013) have greatly broadened our understanding of single
341 gene functions. However, combining the functional annotations of genes to infer the functions of the
342 complex remains challenging and there is no such work published, to the best of our knowledge. In this
343 study, we describe PCfun, the first hybrid computational framework for function prediction of protein
344 complex, powered by the integration of machine-learning and text-mining techniques. PCfun was built
345 based on large-scale publications obtained from PubMed Central (approximately one million full-text
346 articles and their abstracts). The resulting word-embedding matrix was then used to build two
347 complementary computational models, a supervised RF model for the prediction of function based on
348 the annotations in the CORUM database, and an unsupervised k -d tree model for nearest-neighbor

349 queries. The adapted protein complex leave-one-out cross-validation demonstrated the accurate
350 prediction performance of the supervised RF model for predicting the PC-GO associations with respect
351 to the ground-truth annotations provided in CORUM. On the other hand, we also constructed the k -d tree
352 model, which is able to shortlist top-ranked nearest-neighbors, including terms not annotated in the
353 CORUM database. Through enrichment analysis of the functional terms predicted by both models,
354 PCfun is able to provide the final predicted GO terms associated with the input protein complex. Its
355 representation of a contextual GO dendrogram structure along with the embedded predicted GO terms
356 illustrates the hierarchal relationships of all predicted GO terms for the given protein complex.

357 As discussed in the ‘Results’ section, one issue during the construction of the machine-learning
358 models is the biased functional annotations in the CORUM database, as shown in **Figure S3**. Therefore,
359 the predictive power of the RF model in PCfun is limited to the CORUM annotations, demonstrating
360 that it is crucial to combine the ‘non-biased’ prediction results of unsupervised k -d tree method with RF
361 predictions via enrichment analysis, in order for PCfun to deliver non-biased predicted functions for a
362 given protein complex. Another noteworthy issue is the negative data for training the supervised RF
363 model of PCfun. Traditional supervised-learning models for binary classification require accurate
364 labeling of classes (e.g. positive vs. negative) for each sample in the training dataset. However, labelling
365 negative training samples is practically challenging in biological context due to the lack of experimental
366 data. Oftentimes a previously labelled negative sample can be relabeled as positive with the acquisition
367 of new data from novel biological techniques or identification methods. We will re-train the PCfun model
368 once the annotations in the referred databases are updated, in order to continually improve the prediction
369 performance.

370 PCfun can be applied to broad biological and personalized medicine applications. As
371 demonstrated in this work, PCfun is easily accessible and offers utility to a variety of proteomic
372 technologies. Prior to the prediction, the potential protein complexes with clearly defined subunits (either
373 UniProt accessions or gene names) can be provided using the methodology described in the study of
374 McBride *et al* (McBride et al., 2019). It is also possible to compare the functional differences by
375 examining the prediction outputs of PCfun for a protein complex across different biological/medical

376 conditions. In the future, we also plan to incorporate the differential analysis, similar to gene set
377 enrichment analysis, for differentially regulated protein complexes that may have the ability to leverage
378 both compositional (i.e. stoichiometric) and abundance differences. Additionally, prior to delivering the
379 prediction results, PCfun first searches the given protein complex within the CORUM database for
380 possible annotations. The documented annotations in CORUM and newly predicted GO terms are
381 separated in the final outputs in an effort to facilitate the generation of novel hypotheses regarding the
382 function of the protein complex. Taken together, we anticipate that PCfun can serve as an instrumental
383 computational approach for the prediction of novel functions of protein complexes to provide reliable
384 computational evidence for further experimental validation.

385

386 **ACKNOWLEDGEMENTS**

387 This work was supported by in part by the Swiss National Science Foundation (grant No. 3100A0-688
388 107679 to R.A.) and the European Research Council (ERC-2014AdG 670821 to R.A.). C.L. is currently
389 supported by a National Health and Medicine Research Council of Australia (NHMRC) CJ Martin Early
390 Career Research Fellowship (1143366). M.B. was funded by an SNSF SystemsX.ch fellowship (TPdF
391 2013/135). A.W.P. is supported by an NHMRC Principal Research Fellowship (1137739). We thank
392 Fabian Frommelt, Federico Uliana, Michiel Karrenbelt, Elias Pratschke and Shreyans Jain from ETH
393 Zürich for their critical comments and insightful suggestions.

394

395 **AUTHOR CONTRIBUTIONS**

396 R.A., C.L., M.M., V.S.S. and J.S. conceived and designed the project. V.S.S., C.L. and M.M. developed
397 and implemented the PCfun package, and conducted data analysis, machine-learning prediction, and GO
398 term enrichment analysis. M.M. assist with the generation of word-embedding and k -d tree. A.F. and
399 M.M. assisted with package implementation and data analysis. R.C., E.G.W., M.B., W.S., P.G.A.P,
400 M.R.M., Z.C. and A.W.P. provided critical and insightful comments during the development of PCfun.
401 V.S.S., C.L., J.S. and R.A. drafted the manuscript, which has been revised and approved by all the other
402 authors.

403

404 **DECLARATION OF INTERESTS**

405 The authors declare no competing financial interest.

406 STAR METHODS

407 KEY RESOURCES TABLE

REAGENT OR RESOURCE	SOURCE	IDENTIFIER
Deposited Data		
CORUM Database	(Giurgiu et al., 2019)	https://mips.helmholtz-muenchen.de/corum/
Gene Ontology Resource	(Ashburner et al., 2000; The Gene Ontology, 2019)	http://geneontology.org/
QuickGO database	(Binns et al., 2009)	https://www.ebi.ac.uk/QuickGO/
Complex Portal	(Meldal et al., 2019)	https://www.ebi.ac.uk/complexportal/home
PubMed Central	-	https://www.ncbi.nlm.nih.gov/pmc/
UniProt Database	(UniProt, 2019)	https://www.uniprot.org/
Software and Algorithms		
Scikit-learn	(Pedregosa et al., 2011)	https://scikit-learn.org/stable/
fastText	(Joulin et al., 2017)	https://fasttext.cc/
PCfun	This work	https://github.com/sharmavaruns/PCfun

408

409 LEAD CONTACT AND MATERIALS AVAILABILITY

410 Further information and requests for resources should be directed to and will be fulfilled by the Lead
411 Contact, Chen Li (Chen.Li@monash.edu).

412

413 METHOD DETAILS

414 Text corpus generation and data processing

415 Approximately 1 million articles (including open-access full-text articles and their abstracts) were
416 downloaded from PubMed Central in February 2018. Note that these publications are not
417 species/organism specific, which means that the developed PCfun, built on the corpus, is a generic tool
418 for protein complex function prediction. For processing the articles into a text corpus, we followed the
419 text processing pipeline described in the study of Manica *et al* (Manica et al., 2019). All of the natural
420 language queries were pre-processed by removing all punctuation characters, by fixing Unicode
421 mojibake and garbled HTML entities, and by converting all uppercase characters into lowercase. For the
422 extraction of a single word vector from a natural language query (e.g. a protein complex or GO term
423 name), L2 normalized vectors of the query individual words were extracted from the embedding and
424 averaged. The final averaged vector of the component vectors of the name were L2 normalized again
425 and subsequently used as the final word vector for the natural language query.

426

427 **Word-embedding and similarity calculations**

428 Word-embedding refers to a class of approaches developed in the fields of text mining and natural
429 language processing that embed natural language texts into high-dimensional, continuous real-valued
430 vector representations (Manica et al., 2019). The word-embedding in this study was achieved by training
431 the unsupervised version of ‘fastText’ (Joulin et al., 2017) with a skip-gram model on the publications
432 corpus. Briefly, the skip-gram model attempts to minimize the negative sum of the log probability that a
433 word exists within the context of a target word. The fastText unsupervised training parameters used were
434 the default values as declared in the fastText package except for 500-dimensions for the embedding layer,
435 a context window of size 9, and the usage of bi-grams as chosen from the study of Manica *et al* (Manica
436 et al., 2019). After training, the word vectors were normalized to a unitary norm. We utilized the cosine
437 similarity between two vectors, which measures the cosine of the angle between the two vectors as the
438 normalized dot product of the two vectors.

439

440 **Calculating sub-embeddings of topics and nearest neighbors**

441 We extracted sub-embeddings (a data frame consisting of the word vectors for a particular class of terms,
442 such as all biological process GO terms) for protein complex names (extracted from CORUM- one sub-
443 embedding for each naming scheme: *canonical* or *subunit* names) and GO terms (extracted from GO
444 resource split into biological process, molecular function and cellular component classes). For example,
445 to obtain the GO terms most similar to a protein complex query in the embedding space, we calculated
446 of the nearest neighbors to the protein complex query vector within the extracted GO term sub-
447 embedding space. A nearest-neighbor calculation involves calculating distances between vectors within
448 the sub-embedding of interest and the query vector, and then sorting the neighbors by their similarity by
449 descending order. The calculation for nearest neighbors can be quite time intensive depending on the
450 size of the embedding due to the time requirements for pairwise calculation of the distances between
451 each vector in the sub-embedding and the query vector. Therefore, we stored the sub-embedding vectors
452 into a pre-computed k -d tree (k -dimensional tree), which is a space-partitioning data structure that stores

453 the vectors into buckets determined by hyperplane splits over each dimension of the vector (Freidman et
454 al., 1977). Therefore, calculation of n -nearest neighbors to a query vector requires only placement of the
455 query vector into its corresponding location within the pre-calculated k -d tree and then subsequent
456 querying of its ancestors in the tree until the closest n -nearest neighbors have been calculated.

457

458 **Databases for protein complex annotations**

459 For this work we employed the CORUM database (Giurgiu et al., 2019) as the main resource for the
460 ground-truth annotation of protein complexes with GO terms, as CORUM is a compendium of manually
461 curated and experimentally validated protein complexes for various organisms (Giurgiu et al., 2019).
462 Annotations in CORUM for the function of protein complexes have been collected from various types
463 of evidence, including experimental evidence ('exp'), evidence from literature ('lit'), known mammalian
464 homologs ('kmh'), high-throughput experiments ('htp'), and predicted function ('pred'). Here the
465 'predicted function' refers to the potential function suggested by the experimental results. In this work,
466 we utilized annotations from all species in the CORUM database to keep as much information as possible
467 for constructing an accurate supervised machine-learning model, given that the corpus we obtained from
468 PubMed are not species/organism specific. In our study, non-redundant 3414 core protein complexes
469 (downloaded in March 2019) from the CORUM database were used.

470 In addition to CORUM, we utilized the annotated *Homo sapiens* protein complexes (downloaded
471 in November 2019) from the Complex Portal database (Meldal et al., 2019) to independently assess the
472 prediction performance of PCfun. Similar to CORUM, the Complex Portal contains the protein
473 complexes and their annotations of GO terms. As the Complex Portal has fewer number of protein
474 complexes documented, we did not use it for training the model. For the independent test, only the protein
475 complexes annotated as 'physical interaction evidence used in manual assertion' with the evidence code
476 'ECO:0000353' coupled with experimental evidence from the IntAct database (Orchard et al., 2014)
477 were retained. To objectively benchmark the performance of PCfun on the Complex Portal, we further
478 removed those protein complexes from the Complex Portal that had a subunit overlap of larger than 50%
479 compared to the complexes in the CORUM database. As a result, the numbers of protein complexes from

480 the Complex Portal for the independent test were in total 34, of which 34, 31, and 33 protein complexes
481 have biological process, molecular function and cellular component annotations, respectively.

482

483 **Gene name extraction for protein complex subunits**

484 As each unique natural language query has its own unique word vector, it is important to standardize the
485 natural language name used for each protein complex when extracting its word vector. A protein complex
486 can be either represented as its documented name, as written into the protein complex database, or as its
487 subunits' gene names strung together. In this study, we tested the performance of the algorithms using
488 two different naming schemes: (1) *canonical name* (as documented in CORUM), and (2) *subunit name*
489 (composed by UniProt gene names of each subunit). To obtain the gene names of the subunits, we
490 extracted the subunits of the protein complexes from CORUM and queried the UniProt database
491 (downloaded in May 2019) (UniProt, 2019) for their corresponding gene names using the organism
492 identifier from the CORUM database. We then represented the protein complex name with its text pre-
493 processed subunits' gene names, strung together with spaces demarcating each individual gene name.
494 Due to the fact that there might exist multiple names for a single gene, only its canonical name was
495 extracted and used in our model. Gene names were extracted by downloading the UniProt FASTA-
496 formatted sequence file, with respect to the appropriate species, which was then subsequently parsed for
497 each relevant UniProt ID - gene name pair. The resulting work has tested the performance of the
498 algorithm when using each naming scheme independently from each other. Importantly, use of the
499 subunit UniProt gene name scheme also allows for greater flexibility as one can still gain functional
500 insight into a newly detected protein complex even if the complex has not been officially named yet.

501

502 **Preparation of protein complex - GO (PC-GO) pair datasets for supervised learning**

503 To enable the accurate prediction of protein complex function annotations, we have formulated our task
504 as a supervised binary classification problem, for which we created a labelled dataset based on the
505 CORUM annotations. To create the labelled dataset, we first extracted PC-GO term pairs and then
506 labelled each pair as positive if its annotation was observed in the CORUM database. The GO terms and

507 their DAG structures were collected from the Gene Ontology Resources platform (Ashburner et al., 2000;
508 The Gene Ontology, 2019). To label the negatives, we first generated a pool of all possible negative PC-
509 GO pairs to sample from by taking the GO terms (split by biological process, molecular function and
510 cellular component categories) that were used in CORUM and not annotated for a particular protein
511 complex. Since only a few GO terms were annotated for each protein complex, this negative sample pool
512 significantly outnumbered the positive sample pool. This issue of having a huge number of negatives
513 compared to positive labels is a common problem shared by a variety of bioinformatics studies and have
514 been discussed in the recent studies of Li *et al* (Li et al., 2018b; Li et al., 2019a). To build an unbiased
515 supervised classifier, it is common practice to train on an approximately equal distribution of positive
516 and negative labels. To ensure this, we randomly selected an approximately equal number of negative
517 PC-GO pairs from the pool of negatives as there were positives for each protein complex. This random
518 selection was repeated five times, resulting in five training datasets, where only negative samples were
519 different.

520

521 **Supervised machine-learning classifiers**

522 Supervised binary classification was performed with the RF (Breiman, 2001), LR (Lecessie and
523 Vanhouwelingen, 1992; Yu et al., 2011), and NB (with Gaussian and Bernoulli distributions,
524 respectively) (Zhang, 2004) classifiers with the default parameters. The feature space consisted of 1,000-
525 dimensional vectors (500-dimensional protein complex vector prepended to a 500-dimensional GO term
526 vector) with each vector corresponding to a PC-GO pair labelled with either positive or negative. *RF*
527 classification uses the ensemble of decision trees that randomly bootstraps over the training data and
528 features for each decision tree in order to classify the input vector space. After construction of the random
529 decision trees, the classifier outputs a class membership (positive or negative for binary classification)
530 prediction dependent on the community-wide vote from the random trees constructed. *LR* utilizes the
531 logistic function for the binary classification of a dependent variable. This method outputs a probability
532 score ranging from 0 to 1 where values of >0.5 are considered to be of the positive class membership.
533 *NB classifiers* are probabilistic classifiers built upon Bayesian statistics. These classifiers attempt to learn

534 the distributional fit for each labelled class and accordingly assign probability values for each class when
535 given a new input vector. The priors we tested were the Gaussian and Bernoulli distributions. Briefly,
536 the Gaussian distribution assumes that the continuous data values are distributed according to a normal
537 distribution, whereas the Bernoulli distribution assumes that the features are independent binary
538 variables and can be considered as a special case of the binomial distribution. These two classifiers are
539 termed NB_Gauss and NB_Bernoulli in the following sections, respectively. For PCfun's machine
540 learning classifier predicted terms we used a majority voting scheme over the five datasets to provide an
541 averaged predicted probability for each GO term. For this majority vote, we equally weighted the
542 contribution of each dataset's model to achieve a final combined probability score for particular GO
543 terms. If the combined probability score was >0.5 , the GO term was classified as a positive term,
544 otherwise classified as a negative.

545

546 **Performance assessment of supervised binary classification**

547 To assess the predictive performance of the supervised binary classifiers, we introduced an adaptation
548 on traditional leave-one-out cross-validation, which we termed 'protein complex'-leave-one-out cross
549 validation. The 'protein complex'-leave-one-out cross validation first pre-removed every row that had
550 the particular protein complex being tested (including both positive and negatively labelled rows).
551 Afterward, the model was then trained on the remaining dataset and tested on the pre-removed protein
552 complex of interest's rows. This complex-wise evaluation strategy was applied to each protein complex
553 in the dataset. To measure the performance, we used widely established performance metrics used in a
554 variety of bioinformatics and computational biology studies (Li et al., 2018a; Manavalan et al., 2019;
555 Song et al., 2018), including accuracy, AUC (Area Under the Curve), precision, recall, MCC (Matthews
556 Correlation Coefficient) (Matthews, 1975), and F1 score. All metrics reported in this study are the
557 average of five datasets of 'protein complex'-leave-one-out cross-validation. For plotting the Receiver
558 Operating Characteristic (ROC) curves, we followed the recommendations of the interpolation scikit-
559 learn package (please refer to the scikit-learn tutorial). These scores are calculated from the elementary

560 scores of true positives (*TP*), true negatives (*TN*), false positives (*FP*), and false negatives (*FN*). The
561 formulas for each performance metric are shown as follows:

$$562 \quad Accuracy = \frac{TP + TN}{TP + FP + FN + TN},$$

$$563 \quad Precision = \frac{TP}{TP + FP},$$

$$564 \quad Recall = \frac{TP}{TP + FN},$$

$$565 \quad MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FN)(TP + FP)(TN + FP)(TN + FN)}},$$

$$566 \quad F1 = 2 \times \frac{precision \times recall}{precision + recall}$$

567

568 **Semantic similarity calculation between GO terms**

569 To compare the GO terms predicted by PCfun and the annotations from Complex Portal, we did not
570 perform the direct comparison/matching of the GO terms, due to the low number of intersecting GO
571 terms between CORUM and Complex Portal. Instead, we exploited the semantic similarity introduced
572 by Wang *et al* (Wang et al., 2007). This method considers both biological meanings and hierarchical
573 relationships in the GO DAG structure of the given GO term pair. Any GO term pairs with a semantic
574 similarity ≥ 0.5 were considered similar. When comparing the predicted GO terms by PCfun and the
575 annotations from Complex Portal, for each protein complex, we first calculated the similarities for all
576 possible GO term pairs between the PCfun outputs and Complex Portal annotations. The pairs with
577 similarity score ≥ 0.5 were selected and the unique Complex Portal GO terms (*N*) in the pairs counted.
578 Then the percent coverage of the predicted GO terms to the Complex Portal annotation for the particular
579 protein complex was calculated using $Coverage = (N/M) \times 100$, where *M* denotes the number of GO
580 terms annotated in Complex Portal for the particular protein complex.

581

582

583 **GO term enrichment analysis of predicted functional terms**

584 This step aims to systematically and comprehensively combine the two predicted GO term lists in each
585 category (i.e. biological process, molecular function, and cellular component) by RF and the k -d tree for
586 a given protein complex, respectively, due to the fact that this shortlist of terms predicted by RF that
587 recovers CORUM database well but tends to predict broader GO terms for a protein complex. Given the
588 predicted GO term list by RF with the size N ($N \leq 10$), we supplemented the information from this list
589 with the list of GO terms (i.e. the nearest neighbors) obtained directly by querying the GO term sub-
590 embedding. To accomplish this, we developed a functional enrichment analysis pipeline based on the
591 hypergeometric test to assess if all the child nodes of the GO term by RF are significantly enriched in
592 the predicted GO terms by the k -d tree, using the following formula:

$$593 \quad p = f_{\text{hypergeometric}}(x - 1, M, n, N)$$

594 where M denotes the number of total GO terms for a particular GO term class, n denotes the number of
595 child terms of a parent GO term plus the parent term predicted by the supervised classifier that exists
596 within the specific GO term class, where N denotes the sample size which is the mean number of GO
597 terms required for the nearest neighbor list to recover all GO terms annotated for a protein complex
598 (biological process = 11044, molecular function = 5213, and cellular component = 1896, respectively),
599 and x denotes the number of child terms for a particular supervised term that exist in the set of the sample
600 size list. The function $f_{\text{hypergeometric}}$ was the survival function for the hypergeometric distribution as
601 implemented in the SciPy package. If the child nodes/terms are significantly enriched in the predicted
602 list by the k -d tree, 10 top ranked terms based on cosine similarity from the k -d tree list are selected. We
603 could therefore obtain a ‘combined’ predicted list that not only accurately recovers CORUM database
604 but also supplements the list by RF using the detailed GO terms predicted by the k -d tree. To visualize
605 the results, PCfun plots a GO tree structure of predicted GO term by RF (in green) and the 10 top ranked
606 GO terms by the k -d tree (in purple) to demonstrate the hierarchical relationships of these terms. For the
607 cellular component category, in addition to the combination of the two lists from the k -d tree and the RF
608 model by functional enrichment analysis, we also considered adding the overlap of cellular component

609 annotations of all the subunits to the final outputs. The cellular component annotations for each subunit
610 was downloaded from the QuickGO database (Binns et al., 2009).

611

612 **PCfun prediction output organization**

613 In total, there are six output lists (two for each GO category) for a given protein complex generated by
614 PCfun. For each GO category (i.e. biological process, molecular function, and cellular component), one
615 list contained the RF predictions and the top 10 significantly enriched GO terms by the k -d tree, while
616 the other lists provided the top 20 GO terms by the k -d tree only. In addition, for each RF predicted term,
617 a GO DAG structure is plotted to illustrate the hierarchical relationships between the RF prediction and
618 the top 10 significantly enriched terms from the k -d tree.

619

620 **DATA AND CODE AVAILABILITY**

621 **Data and code availability statement**

622 *Data availability*

623 The full-text articles and their abstracts (in the non-commercial use collection) were extracted from the
624 PubMed Central under a Creative Commons or similar license. The training dataset and the validation
625 test (i.e. the PC-GO association) were obtained from the CORUM database (Giurgiu et al., 2019) and
626 the Complex Portal database (Meldal et al., 2019), respectively. The full GO lists were downloaded from
627 the Gene Ontology Resource platform (Ashburner et al., 2000; The Gene Ontology, 2019). For
628 evaluation purposes, we downloaded the GO annotations for individual proteins from the QuickGO
629 database (Binns et al., 2009).

630

631 *Code availability*

632 PCfun is an open-access software and is freely available for academic use under the ‘Academic Free
633 License v3.0’. The source code, user instruction, and example inputs can be downloaded from
634 <https://github.com/sharmavaruns/PCfun>.

635

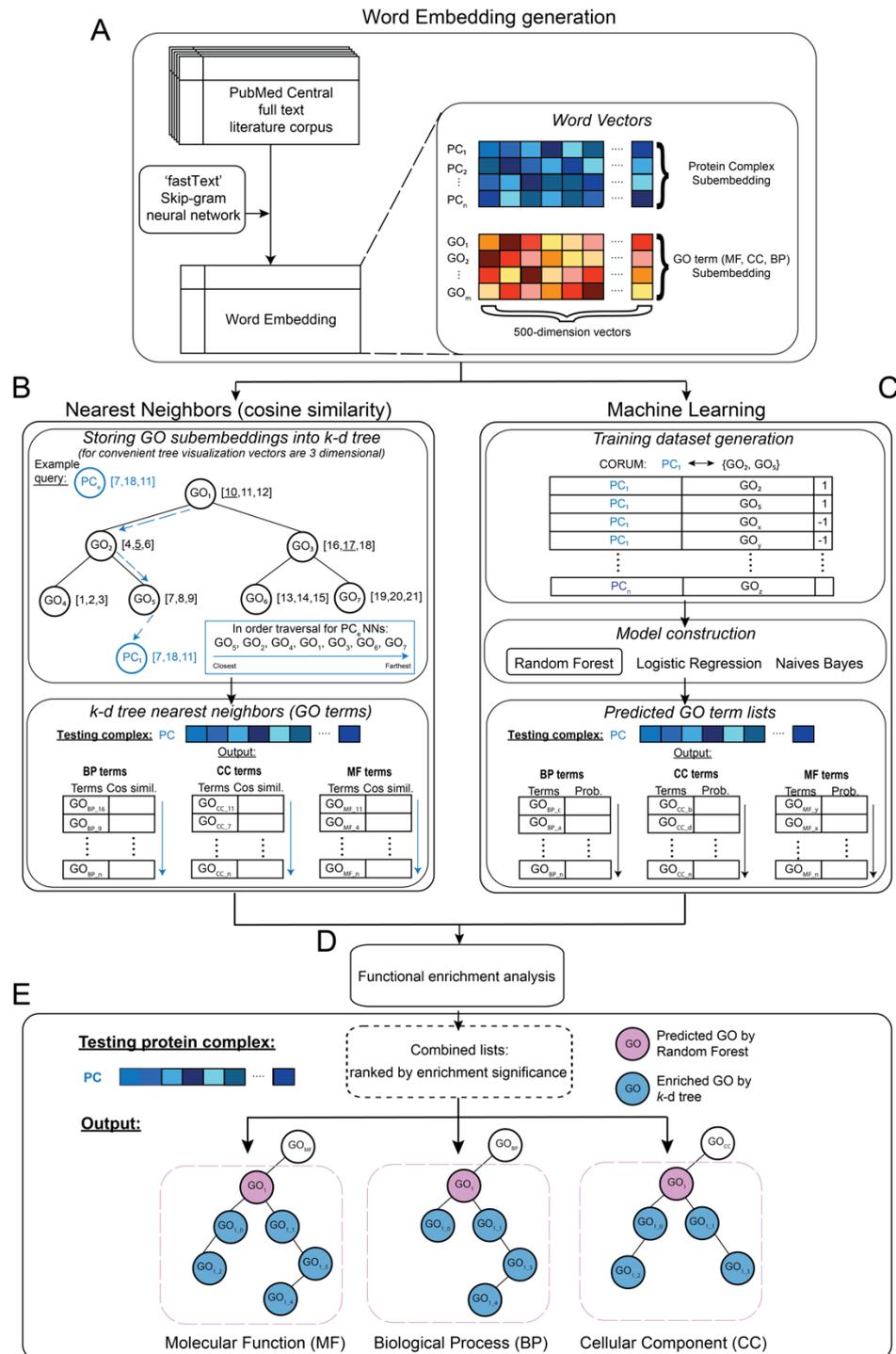
636 **REFERENCES**

- 637 Aebersold, R., and Mann, M. (2016). Mass-spectrometric exploration of proteome structure and function.
638 *Nature* 537, 347-355.
- 639 Aranda, S., Mas, G., and Di Croce, L. (2015). Regulation of gene transcription by Polycomb proteins.
640 *Sci Adv* 1, e1500737.
- 641 Ashburner, M., Ball, C.A., Blake, J.A., Botstein, D., Butler, H., Cherry, J.M., Davis, A.P., Dolinski, K.,
642 Dwight, S.S., Eppig, J.T., *et al.* (2000). Gene ontology: tool for the unification of biology. The Gene
643 Ontology Consortium. *Nat Genet* 25, 25-29.
- 644 Barutcuoglu, Z., Schapire, R.E., and Troyanskaya, O.G. (2006). Hierarchical multi-label prediction of
645 gene function. *Bioinformatics* 22, 830-836.
- 646 Becher, I., Andres-Pons, A., Romanov, N., Stein, F., Schramm, M., Baudin, F., Helm, D., Kurzawa, N.,
647 Mateus, A., Mackmull, M.T., *et al.* (2018). Pervasive Protein Thermal Stability Variation during the Cell
648 Cycle. *Cell* 173, 1495-1507 e1418.
- 649 Binns, D., Dimmer, E., Huntley, R., Barrell, D., O'Donovan, C., and Apweiler, R. (2009). QuickGO: a
650 web-based tool for Gene Ontology searching. *Bioinformatics* 25, 3045-3046.
- 651 Breiman, L. (2001). Random forests. *Mach Learn* 45, 5-32.
- 652 Carlson, M.L., Stacey, R.G., Young, J.W., Wason, I.S., Zhao, Z., Rattray, D.G., Scott, N., Kerr, C.H.,
653 Babu, M., Foster, L.J., *et al.* (2019). Profiling the Escherichia coli membrane protein interactome
654 captured in Peptidisc libraries. *Elife* 8.
- 655 Cejuela, J.M., Vinchurkar, S., Goldberg, T., Prabhu Shankar, M.S., Baghudana, A., Bojchevski, A.,
656 Uhlig, C., Ofner, A., Raharja-Liu, P., Jensen, L.J., *et al.* (2018). LocText: relation extraction of protein
657 localizations to assist database curation. *BMC Bioinformatics* 19, 15.
- 658 Chen, Y., Chen, S., Li, K., Zhang, Y., Huang, X., Li, T., Wu, S., Wang, Y., Carey, L.B., and Qian, W.
659 (2019). Overdosage of Balanced Protein Complexes Reduces Proliferation Rate in Aneuploid Cells. *Cell*
660 *Syst* 9, 129-142 e125.
- 661 D'Avino, P.P., Archambault, V., Przewloka, M.R., Zhang, W., Laue, E.D., and Glover, D.M. (2009).
662 Isolation of protein complexes involved in mitosis and cytokinesis from Drosophila cultured cells.
663 *Methods Mol Biol* 545, 99-112.
- 664 Freidman, J.H., Bentley, J.L., and Finkel, R.A. (1977). An Algorithm for Finding Best Matches in
665 Logarithmic Expected Time. *ACM Transactions on Mathematical Software* 3, 209.
- 666 Gaizauskas, R., Demetriou, G., Artymiuk, P.J., and Willett, P. (2003). Protein structures and information
667 extraction from biological texts: the PASTA system. *Bioinformatics* 19, 135-143.
- 668 Giurgiu, M., Reinhard, J., Brauner, B., Dunger-Kaltenbach, I., Fobo, G., Frishman, G., Montrone, C.,
669 and Ruepp, A. (2019). CORUM: the comprehensive resource of mammalian protein complexes-2019.
670 *Nucleic Acids Res* 47, D559-D563.
- 671 Guan, Y., Myers, C.L., Hess, D.C., Barutcuoglu, Z., Caudy, A.A., and Troyanskaya, O.G. (2008).
672 Predicting gene function in a hierarchical context with an ensemble of classifiers. *Genome Biol* 9 *Suppl*
673 *1*, S3.
- 674 Hartwell, L.H., Hopfield, J.J., Leibler, S., and Murray, A.W. (1999). From molecular to modular cell
675 biology. *Nature* 402, C47-52.
- 676 Heusel, M., Bludau, I., Rosenberger, G., Hafen, R., Frank, M., Banaei-Esfahani, A., van Drogen, A.,
677 Collins, B.C., Gstaiger, M., and Aebersold, R. (2019). Complex-centric proteome profiling by SEC-
678 SWATH-MS. *Mol Syst Biol* 15, e8438.
- 679 Heusel, M., Frank, M., Kohler, M., Amon, S., Frommelt, F., Rosenberger, G., Bludau, I., Aulakh, S.,
680 Linder, M.I., Liu, Y., *et al.* (2020). A Global Screen for Assembly State Changes of the Mitotic Proteome
681 by SEC-SWATH-MS. *Cell Syst* 10, 133-155 e136.

- 682 Hewick, R.M., Lu, Z., and Wang, J.H. (2003). Proteomics in drug discovery. *Adv Protein Chem* *65*, 309-
683 342.
- 684 Islamaj Dogan, R., Kim, S., Chatr-Aryamontri, A., Wei, C.H., Comeau, D.C., Antunes, R., Matos, S.,
685 Chen, Q., Elangovan, A., Panyam, N.C., *et al.* (2019). Overview of the BioCreative VI Precision
686 Medicine Track: mining protein interactions and mutations for precision medicine. *Database (Oxford)*
687 *2019*.
- 688 Jeffery, C.J. (2015). Why study moonlighting proteins? *Front Genet* *6*, 211.
- 689 Joulin, A., Grave, E., Bojanowski, P., and Mikolov, T. (2017). Bag of Tricks for Efficient Text
690 Classification. Paper presented at: The 15th Conference of the European Chapter of the Association for
691 Computational Linguistics (Valencia, Spain: Association for Computational Linguistics).
- 692 Lecessie, S., and Vanhouwelingen, J.C. (1992). Ridge Estimators in Logistic-Regression. *Appl Stat-J*
693 *Roy St C* *41*, 191-201.
- 694 Leitner, A., Faini, M., Stengel, F., and Aebersold, R. (2016). Crosslinking and Mass Spectrometry: An
695 Integrated Technology to Understand the Structure and Function of Molecular Machines. *Trends*
696 *Biochem Sci* *41*, 20-32.
- 697 Leitner, A., Reischl, R., Walzthoeni, T., Herzog, F., Bohn, S., Forster, F., and Aebersold, R. (2012).
698 Expanding the chemical cross-linking toolbox by the use of multiple proteases and enrichment by size
699 exclusion chromatography. *Mol Cell Proteomics* *11*, M111 014126.
- 700 Li, F., Li, C., Marquez-Lago, T.T., Leier, A., Akutsu, T., Purcell, A.W., Ian Smith, A., Lithgow, T., Daly,
701 R.J., Song, J., *et al.* (2018a). Quokka: a comprehensive tool for rapid and accurate prediction of kinase
702 family-specific phosphorylation sites in the human proteome. *Bioinformatics* *34*, 4223-4231.
- 703 Li, F., Wang, Y., Li, C., Marquez-Lago, T.T., Leier, A., Rawlings, N.D., Haffari, G., Revote, J., Akutsu,
704 T., Chou, K.C., *et al.* (2018b). Twenty years of bioinformatics research for protease-specific substrate
705 and cleavage site prediction: a comprehensive revisit and benchmarking of existing methods. *Brief*
706 *Bioinform.*
- 707 Li, F., Zhang, Y., Purcell, A.W., Webb, G.I., Chou, K.C., Lithgow, T., Li, C., and Song, J. (2019a).
708 Positive-unlabelled learning of glycosylation sites in the human proteome. *BMC Bioinformatics* *20*, 112.
- 709 Li, M., He, Q., Ma, J., He, F., Zhu, Y., Chang, C., and Chen, T. (2019b). PPICurator: A Tool for
710 Extracting Comprehensive Protein-Protein Interaction Information. *Proteomics* *19*, e1800291.
- 711 Liberati, N.T., Datto, M.B., Frederick, J.P., Shen, X., Wong, C., Rougier-Chapman, E.M., and Wang,
712 X.F. (1999). Smads bind directly to the Jun family of AP-1 transcription factors. *Proc Natl Acad Sci U*
713 *S A* *96*, 4844-4849.
- 714 Liu, F., Rijkers, D.T., Post, H., and Heck, A.J. (2015). Proteome-wide profiling of protein assemblies by
715 cross-linking mass spectrometry. *Nat Methods* *12*, 1179-1184.
- 716 Manavalan, B., Basith, S., Shin, T.H., Wei, L., and Lee, G. (2019). mAHTPred: a sequence-based meta-
717 predictor for improving the prediction of anti-hypertensive peptides using effective feature
718 representation. *Bioinformatics* *35*, 2757-2765.
- 719 Manica, M., Mathis, R., Cadow, J., and Martínez, M.R. (2019). Context-specific interaction networks
720 from vector representation of words. *Nature Machine Intelligence* *1*, 10.
- 721 Massague, J., Seoane, J., and Wotton, D. (2005). Smad transcription factors. *Genes Dev* *19*, 2783-2810.
- 722 Matalon, O., Horovitz, A., and Levy, E.D. (2014). Different subunits belonging to the same protein
723 complex often exhibit discordant expression levels and evolutionary properties. *Curr Opin Struct Biol*
724 *26*, 113-120.
- 725 Mateus, A., Kurzawa, N., Becher, I., Sridharan, S., Helm, D., Stein, F., Typas, A., and Savitski, M.M.
726 (2020). Thermal proteome profiling for interrogating protein interactions. *Mol Syst Biol* *16*, e9232.

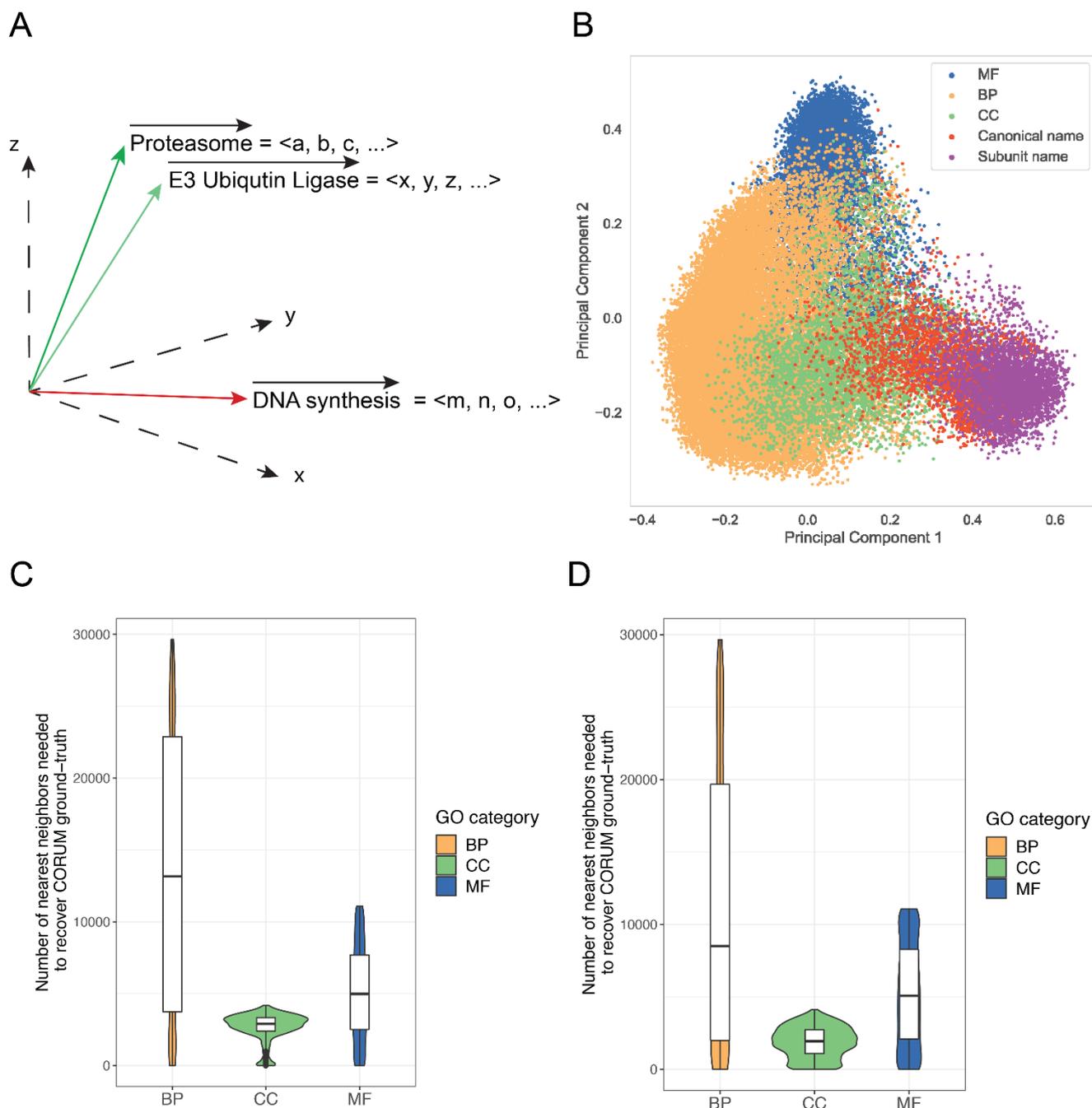
- 727 Matthews, B.W. (1975). Comparison of the predicted and observed secondary structure of T4 phage
728 lysozyme. *Biochim Biophys Acta* *405*, 442-451.
- 729 McBride, Z., Chen, D., Lee, Y., Aryal, U.K., Xie, J., and Szymanski, D.B. (2019). A Label-free Mass
730 Spectrometry Method to Predict Endogenous Protein Complex Composition. *Mol Cell Proteomics* *18*,
731 1588-1606.
- 732 Meldal, B.H.M., Bye, A.J.H., Gajdos, L., Hammerova, Z., Horackova, A., Melicher, F., Perfetto, L.,
733 Pokorny, D., Lopez, M.R., Turkova, A., *et al.* (2019). Complex Portal 2018: extended content and
734 enhanced visualization tools for macromolecular complexes. *Nucleic Acids Res* *47*, D550-D558.
- 735 Nakabayashi, Y., Kawashima, S., Enomoto, T., Seki, M., and Horikoshi, M. (2014). Roles of common
736 subunits within distinct multisubunit complexes. *Proc Natl Acad Sci U S A* *111*, 699-704.
- 737 Orchard, S., Ammari, M., Aranda, B., Breuza, L., Briganti, L., Broackes-Carter, F., Campbell, N.H.,
738 Chavali, G., Chen, C., del-Toro, N., *et al.* (2014). The MIntAct project--IntAct as a common curation
739 platform for 11 molecular interaction databases. *Nucleic Acids Res* *42*, D358-363.
- 740 Pawson, T., and Nash, P. (2000). Protein-protein interactions define specificity in signal transduction.
741 *Genes Dev* *14*, 1027-1047.
- 742 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
743 Prettenhofer, P., Weiss, R., Dubourg, V., *et al.* (2011). Scikit-learn: Machine Learning in Python. *J Mach*
744 *Learn Res* *12*, 2825-2830.
- 745 Pletscher-Frankild, S., Palleja, A., Tsafo, K., Binder, J.X., and Jensen, L.J. (2015). DISEASES: text
746 mining and data integration of disease-gene associations. *Methods* *74*, 83-89.
- 747 Rebois, R.V., and Hebert, T.E. (2003). Protein complexes involved in heptahelical receptor-mediated
748 signal transduction. *Receptors Channels* *9*, 169-194.
- 749 Rosenberger, G., Heusel, M., Bludau, I., Collins, B.C., Martelli, C., Williams, E.G., Xue, P., Liu, Y.,
750 Aebersold, R., and Califano, A. (2020). SECAT: Quantifying Protein Complex Dynamics across Cell
751 States by Network-Centric Analysis of SEC-SWATH-MS Profiles. *Cell Syst* *11*, 589-607 e588.
- 752 Schopper, S., Kahraman, A., Leuenberger, P., Feng, Y., Piazza, I., Muller, O., Boersema, P.J., and Picotti,
753 P. (2017). Measuring protein structural changes on a proteome-wide scale using limited proteolysis-
754 coupled mass spectrometry. *Nat Protoc* *12*, 2391-2410.
- 755 Simonis, N., van Helden, J., Cohen, G.N., and Wodak, S.J. (2004). Transcriptional regulation of protein
756 complexes in yeast. *Genome Biol* *5*, R33.
- 757 Song, J., Li, F., Leier, A., Marquez-Lago, T.T., Akutsu, T., Haffari, G., Chou, K.C., Webb, G.I., Pike,
758 R.N., and Hancock, J. (2018). PROSPEROUS: high-throughput prediction of substrate cleavage sites for
759 90 proteases with improved accuracy. *Bioinformatics* *34*, 684-687.
- 760 Stacey, R.G., Skinnider, M.A., Scott, N.E., and Foster, L.J. (2017). A rapid and accurate approach for
761 prediction of interactomes from co-elution data (PrInCE). *BMC Bioinformatics* *18*, 457.
- 762 Stojanova, D., Ceci, M., Malerba, D., and Dzeroski, S. (2013). Using PPI network autocorrelation in
763 hierarchical multi-label classification trees for gene function prediction. *BMC Bioinformatics* *14*, 285.
- 764 Subramani, S., Kalpana, R., Monickaraj, P.M., and Natarajan, J. (2015). HPIminer: A text mining system
765 for building and visualizing human protein interaction networks and pathways. *J Biomed Inform* *54*,
766 121-131.
- 767 Szklarczyk, D., Gable, A.L., Lyon, D., Junge, A., Wyder, S., Huerta-Cepas, J., Simonovic, M., Doncheva,
768 N.T., Morris, J.H., Bork, P., *et al.* (2019). STRING v11: protein-protein association networks with
769 increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic*
770 *Acids Res* *47*, D607-D613.
- 771 Tan, K., Shlomi, T., Feizi, H., Ideker, T., and Sharan, R. (2007). Transcriptional regulation of protein
772 complexes within and across species. *Proc Natl Acad Sci U S A* *104*, 1283-1288.

- 773 The Gene Ontology, C. (2019). The Gene Ontology Resource: 20 years and still GOing strong. *Nucleic*
774 *Acids Res* *47*, D330-D338.
- 775 UniProt, C. (2019). UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res* *47*, D506-D515.
- 776 Wang, J.Z., Du, Z., Payattakool, R., Yu, P.S., and Chen, C.F. (2007). A new method to measure the
777 semantic similarity of GO terms. *Bioinformatics* *23*, 1274-1281.
- 778 Webb, E.C., and Westhead, D.R. (2009). The transcriptional regulation of protein complexes; a cross-
779 species perspective. *Genomics* *94*, 369-376.
- 780 Wicks, S.J., Lui, S., Abdel-Wahab, N., Mason, R.M., and Chantry, A. (2000). Inactivation of smad-
781 transforming growth factor beta signaling by Ca(2+)-calmodulin-dependent protein kinase II. *Mol Cell*
782 *Biol* *20*, 8103-8111.
- 783 Yu, H.F., Huang, F.L., and Lin, C.J. (2011). Dual coordinate descent methods for logistic regression and
784 maximum entropy models. *Mach Learn* *85*, 41-75.
- 785 Yu, K., Lung, P.Y., Zhao, T., Zhao, P., Tseng, Y.Y., and Zhang, J. (2018). Automatic extraction of
786 protein-protein interactions using grammatical relationship graph. *BMC Med Inform Decis Mak* *18*, 42.
- 787 Zhang, H. (2004). The Optimality of Naïve Bayes. Paper presented at: THE SEVENTEENTH
788 INTERNATIONAL FLORIDA ARTIFICIAL INTELLIGENCE RESEARCH SOCIETY
789 CONFERENCE (AAAI).
- 790



791

792 **Figure 1.** The overall framework of PCfun. (A) A word-embedding containing a 500-dimensional vector
 793 for each word was first generated based on the open-access full-text articles and their abstracts using the
 794 'fastText' with a skip-gram model. Based on the word-embedding, two machine-learning algorithms
 795 were used, including (B) a *k-d* tree for nearest-neighbor search, and (C) a supervised RF model for PC
 796 association with MF, BP and CC, respectively. A simplified *k-d* tree example is shown in the top panel
 797 of (B). To combine the outputs of the two models, GO terms enrichment analysis (D) was performed.
 798 PCfun utilizes the enrichment analysis and GO DAG structure (E) to represent and visualize the predicted
 799 GO terms for a given protein complex. The testing protein complex (i.e. 'PC') is given to illustrate the
 800 usage of PCfun.

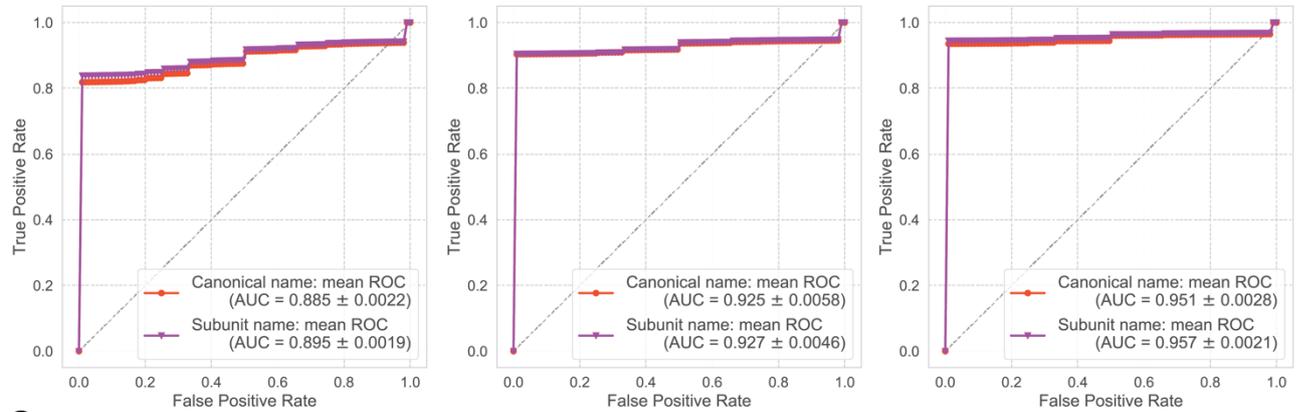


801 For complexes with canonical names For complexes with subunit names
 802
 803 **Figure 2.** Using the word-embedding and k -d tree to shortlist GO terms for protein complexes. (A) A
 804 graphical illustration of the vectors of phrases ‘Proteasome’, ‘E3 Ubiquitin Ligase’ and ‘DNA synthesis’
 805 in a 3D space using simplified vector representation, for example, $\langle a, b, c, \dots \rangle$ denotes the numerical
 806 word vector for the natural language query ‘Proteasome’. (B) The principal component analysis (PCA)
 807 results of different types of sub-embeddings, including molecular function, biological process, cellular
 808 component, protein complexes with canonical names, and subunit names. (C) The numbers of nearest
 809 neighbors required from the k -d tree search outputs for protein complexes using canonical names to cover
 810 the CORUM ground-truth. (D) The numbers of nearest neighbors required from the k -d tree search
 811 outputs for protein complexes using subunit names to cover the CORUM ground-truth.
 812

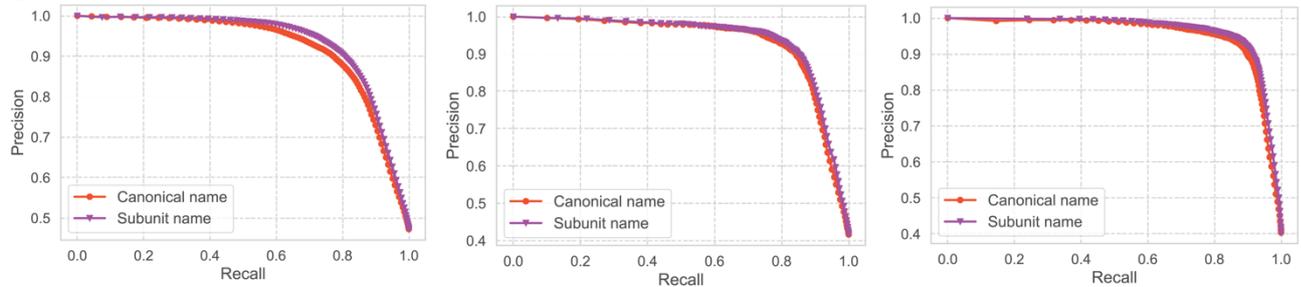
A

Naming scheme	Biological process					Molecular function					Cellular component				
	Accuracy	Precision	Recall	MCC	F1	Accuracy	Precision	Recall	MCC	F1	Accuracy	Precision	Recall	MCC	F1
Canonical name	84.27% ±0.15%	0.748 ±0.002	0.741 ±0.002	0.654 ±0.004	0.729 ±0.002	88.88% ±0.32%	0.799 ±0.004	0.805 ±0.006	0.761 ±0.007	0.795 ±0.004	92.22% ±0.09%	0.862 ±0.002	0.882 ±0.003	0.836 ±0.001	0.866 ±0.002
Subunit name	85.57% ±0.18%	0.769 ±0.003	0.757 ±0.002	0.682 ±0.003	0.748 ±0.002	89.78% ±0.31%	0.812 ±0.006	0.822 ±0.005	0.780 ±0.006	0.811 ±0.005	92.65% ±0.12%	0.868 ±0.001	0.883 ±0.003	0.845 ±0.002	0.870 ±0.001

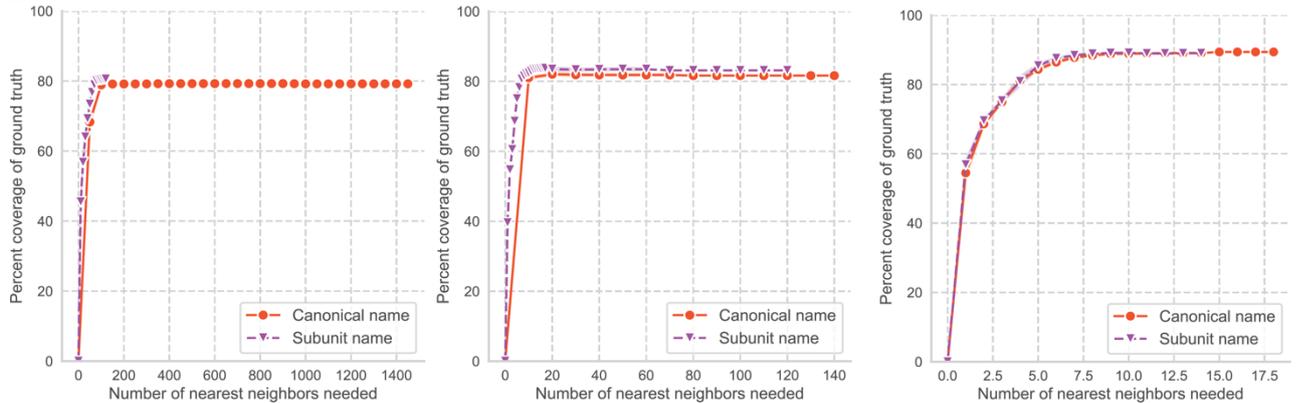
B



C



D



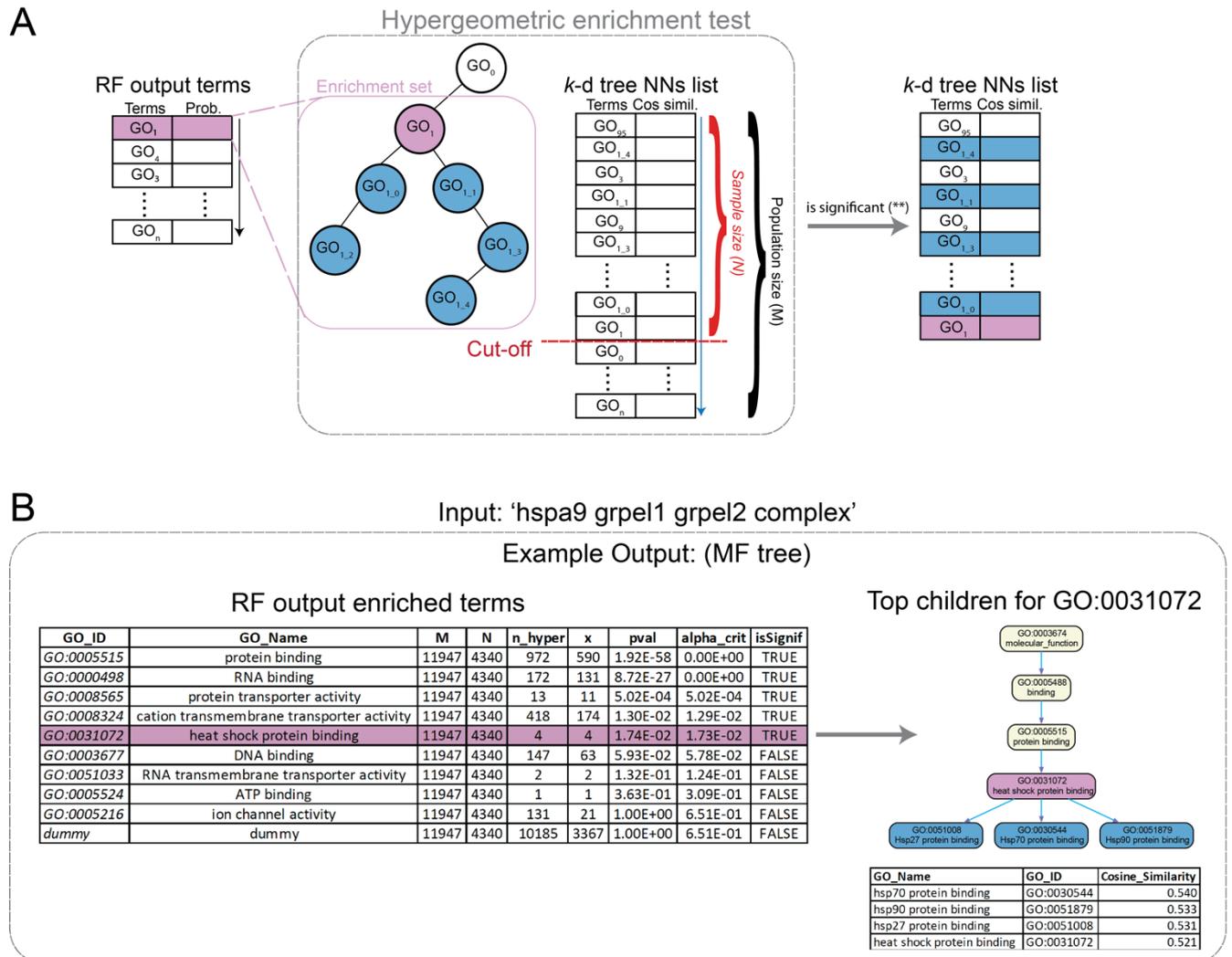
Biological process

Molecular function

Cellular component

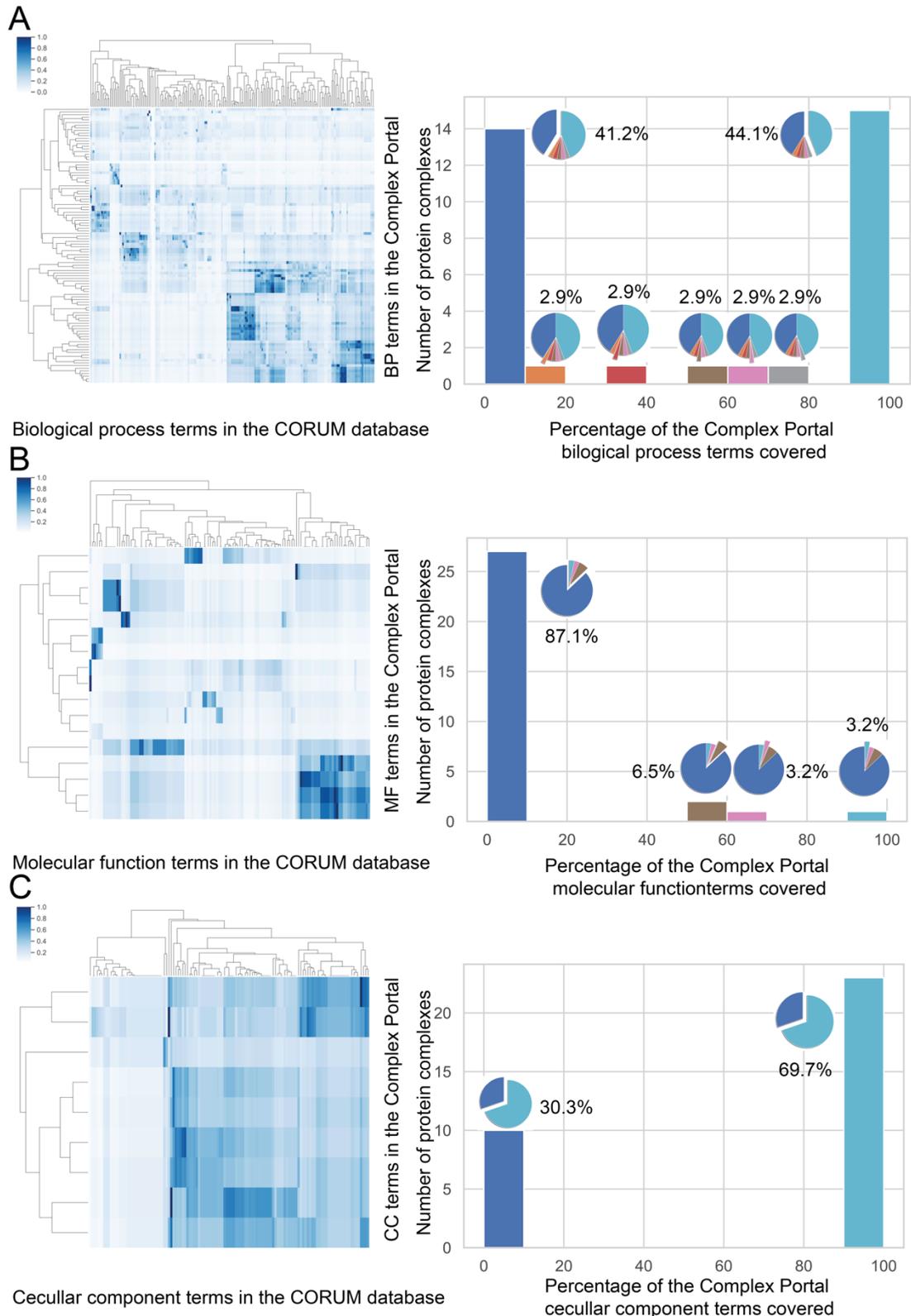
813
814 **Figure 3.** Prediction performance of RF model trained on the ground-truth protein complex – GO term
815 annotations in the CORUM database for biological process, molecular function, and cellular component
816 using canonical and subunit naming schemes for protein complexes, respectively, including (A)
817 performance measures of RF models; (B) ROC curves and mean AUC values; (C) the precision-recall
818 curves of the RF models via the adapted protein complex leave-one-out cross-validation, and (D) the
819 numbers of predicted GO terms by the RF models to recover the CORUM database annotations.

820



821
 822 **Figure 4.** Utilizing functional enrichment analysis to combine the prediction lists from the k -d tree and
 823 the supervised RF model. (A) Hypergeometric enrichment test based on the RF prediction list. For each
 824 term in the list, all the child nodes of the term were collected and used for the statistical significance test
 825 with the terms from the k -d tree. If significant, the top ten terms from the k -d tree that are the child nodes
 826 of the RF term were selected and visualized in the GO DAG structure. (B) An example of prediction
 827 results for the protein complex 'hspa9 grpel1 grpel2 complex' illustrating the enrichment analysis
 828 procedure.

829
 830



831
 832 **Figure 5.** Semantic similarity comparison of (A) biological process, (B) molecular function, and (C)
 833 cellular component terms for the Complex Portal annotations with the CORUM annotations on identical
 834 protein complexes (left panel) and PCfun predictions on non-overlapping protein complexes (right panel).
 835 The right panel illustrates the percentage breakdowns of the predicted similar biological process and
 836 molecular function terms by PCfun to the Complex Portal annotations, respectively. Any GO term pairs
 837 between PCfun predictions and Complex Portal annotations with the semantic similarity ≥ 0.5 were
 838 considered similar.

839 **Supplementary information**

840

841 **Table S1.** Top 10 GO terms shortlisted by the *k*-d tree for the protein complex “SMAD2-SMAD4-
842 FAST1-TGIF-HDAC1 complex, TGF(beta) induced”.

Number	GO term	Cosine distance	Cosine similarity	GO ID
Biological process				
1	jun phosphorylation	0.90142032	0.52592264	GO:0007258
2	common partner smad protein phosphorylation	0.91515871	0.52214994	GO:0007182
3	smad protein signal transduction	0.9206383	0.52066024	GO:0060395
4	pathway restricted smad protein phosphorylation	0.92724758	0.51887469	GO:0060389
5	regulation of histone h3 t3 phosphorylation	0.93482421	0.51684282	GO:2000281
6	regulation of smad protein signal transduction	0.93509793	0.51676971	GO:0007184
7	histone h3 t3 phosphorylation	0.93731143	0.51617927	GO:0072355
8	negative regulation of smad protein signal transduction	0.94274921	0.51473448	GO:0060392
9	regulation of pathway restricted smad protein phosphorylation	0.9430544	0.51465363	GO:0060393
10	negative regulation of histone h3 k27 trimethylation	0.94409184	0.51437899	GO:1902465
Molecular function				
1	smad binding	0.81312028	0.55153539	GO:0046332
2	i smad binding	0.88015558	0.53187088	GO:0070411
3	r smad binding	0.89042146	0.52898257	GO:0070412
4	co smad binding	0.89897215	0.52660067	GO:0070410
5	bmp (bone morphogenic protein) binding	0.93230116	0.51751767	GO:0036122
6	activin binding	0.95586228	0.51128344	GO:0048185
7	histone deacetylase activity h3 k14 specific	0.96414003	0.50912867	GO:0031078
8	bmp receptor binding	0.96599249	0.50864894	GO:0070700
9	transcription corepressor binding	0.97184839	0.50713838	GO:0001222
10	bmp receptor activity	0.97334295	0.50675429	GO:0098821
Cellular component				
1	pdx1 pbx1b mrg1 complex	0.84789228	0.54115709	GO:0034978
2	gata2 tal1 tcf3 lmo2 complex	0.86410241	0.53645121	GO:0070354
3	rgs6 dnmt1 dmap1 complex	0.86604966	0.53589142	GO:0070313
4	gata1 tal1 tcf3 lmo2 complex	0.87113766	0.53443422	GO:0070353
5	heteromeric smad protein complex	0.8795976	0.53202877	GO:0071145
6	smad protein complex	0.88390432	0.53081252	GO:0071141
7	maml3 rbp jkappa icn1 complex	0.89328291	0.52818308	GO:0071179
8	maml1 rbp jkappa icn1 complex	0.89389121	0.52801343	GO:0002193
9	homomeric smad protein complex	0.90011405	0.5262842	GO:0071143
10	fhl2 creb complex	0.90109935	0.52601144	GO:0034980

843

844

345

346

Table S2. Performance comparison of the Logistic Regression (LR) and Naïve Bayes (NB_Gauss and NB_Bernoulli) classifiers trained on the CORUM

347

database via the adapted protein-complex leave-one-out cross-validation, for biological process, molecular function and cellular component categories,

348

respectively.

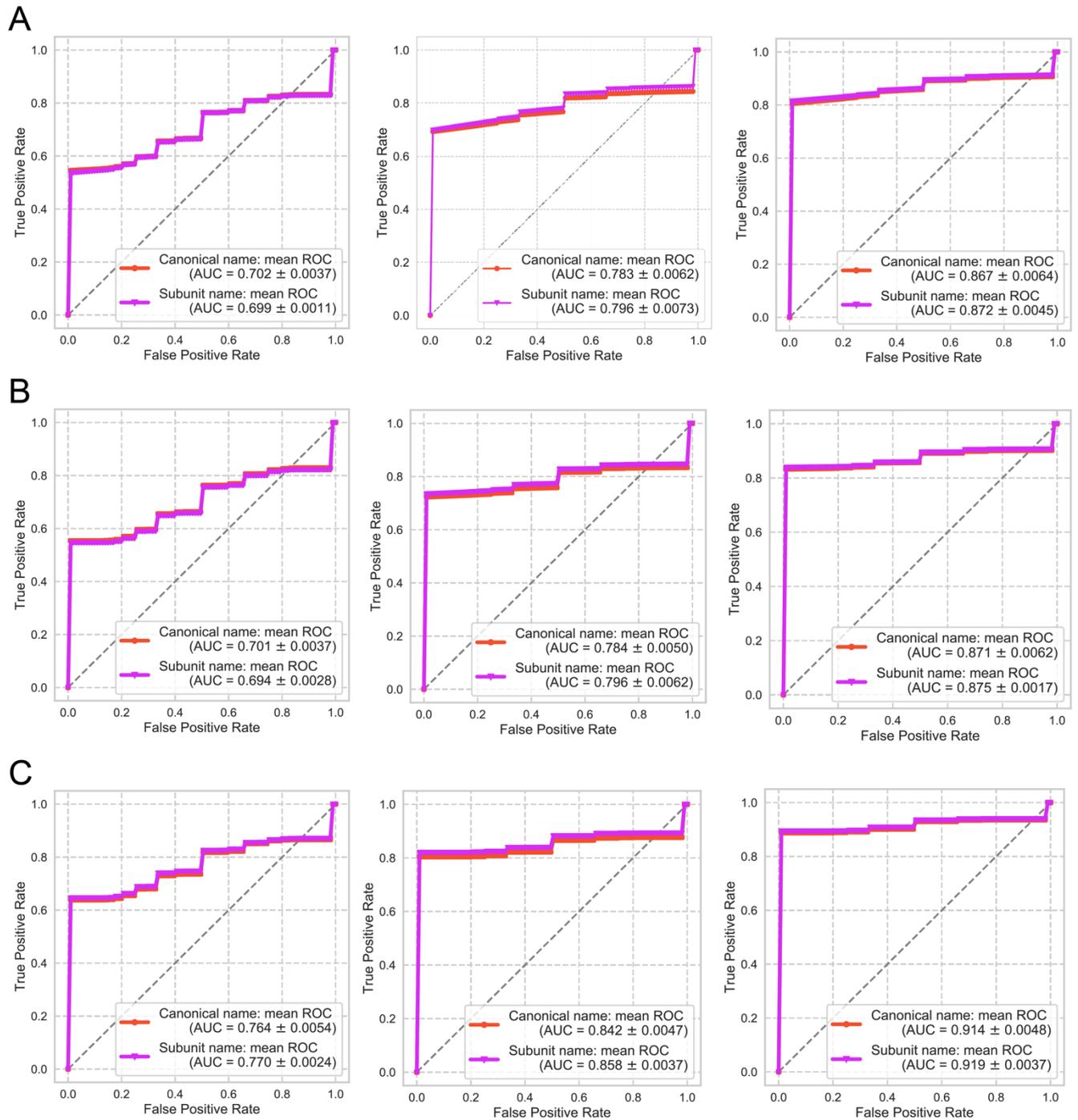
Naming scheme	Biological process						Molecular function						Cellular component					
	Acc. ¹	AUC	Precision	Recall	MCC	F1	Acc.	AUC	Precision	Recall	MCC	F1	Acc.	AUC	Precision	Recall	MCC	F1
LR																		
Canonical name	73.05% ±0.47%	0.764 ±0.006	0.613 ±0.008	0.638 ±0.010	0.419 ±0.009	0.602 ±0.009	81.2% ±0.33%	0.842 ±0.005	0.69 ±0.007	0.708 ±0.006	0.599 ±0.007	0.688 ±0.007	87.55% ±0.34%	0.914 ±0.005	0.785 ±0.005	0.801 ±0.007	0.735 ±0.007	0.785 ±0.006
Subunit name	73.4% ±0.30%	0.77 ±0.003	0.623 ±0.006	0.642 ±0.008	0.428 ±0.004	0.609 ±0.007	82.36% ±0.55%	0.858 ±0.004	0.71 ±0.005	0.727 ±0.005	0.625 ±0.007	0.708 ±0.004	87.71% ±0.40%	0.919 ±0.004	0.789 ±0.007	0.808 ±0.007	0.741 ±0.009	0.79 ±0.008
NB_Gauss																		
Canonical name	64.66% ±0.36%	0.702 ±0.004	0.503 ±0.009	0.547 ±0.012	0.249 ±0.008	0.497 ±0.009	73.89% ±0.375%	0.783 ±0.007	0.594 ±0.014	0.648 ±0.013	0.467 ±0.010	0.602 ±0.012	83.15% ±0.10%	0.867 ±0.007	0.687 ±0.003	0.704 ±0.005	0.638 ±0.002	0.687 ±0.004
Subunit name	64.82% ±0.39%	0.699 ±0.001	0.485 ±0.007	0.529 ±0.009	0.245 ±0.008	0.479 ±0.008	73.54% ±0.672%	0.796 ±0.008	0.581 ±0.012	0.641 ±0.009	0.46 ±0.014	0.593 ±0.010	82.47% ±0.98%	0.872 ±0.005	0.661 ±0.022	0.674 ±0.022	0.619 ±0.022	0.66 ±0.022
NB_Bernoulli																		
Canonical name	64.89% ±0.36%	0.701 ±0.004	0.500 ±0.009	0.534 ±0.008	0.252 ±0.009	0.489 ±0.007	72.56% ±0.658%	0.784 ±0.006	0.563 ±0.013	0.622 ±0.011	0.438 ±0.013	0.575 ±0.011	81.98% ±0.21%	0.871 ±0.007	0.645 ±0.002	0.650± 0.0	0.604 ±0.002	0.642 ±0.001
Subunit name	64.33% ±0.64%	0.694 ±0.003	0.475 ±0.007	0.515 ±0.008	0.236 ±0.012	0.468 ±0.008	72.46% ±0.640%	0.796 ±0.007	0.558 ±0.004	0.621 ±0.0	0.437 ±0.014	0.572 ±0.003	82.18% ±0.21%	0.875 ±0.002	0.647 ±0.001	0.654 ±0.0	0.609 ±0.002	0.644 ±0.001

349

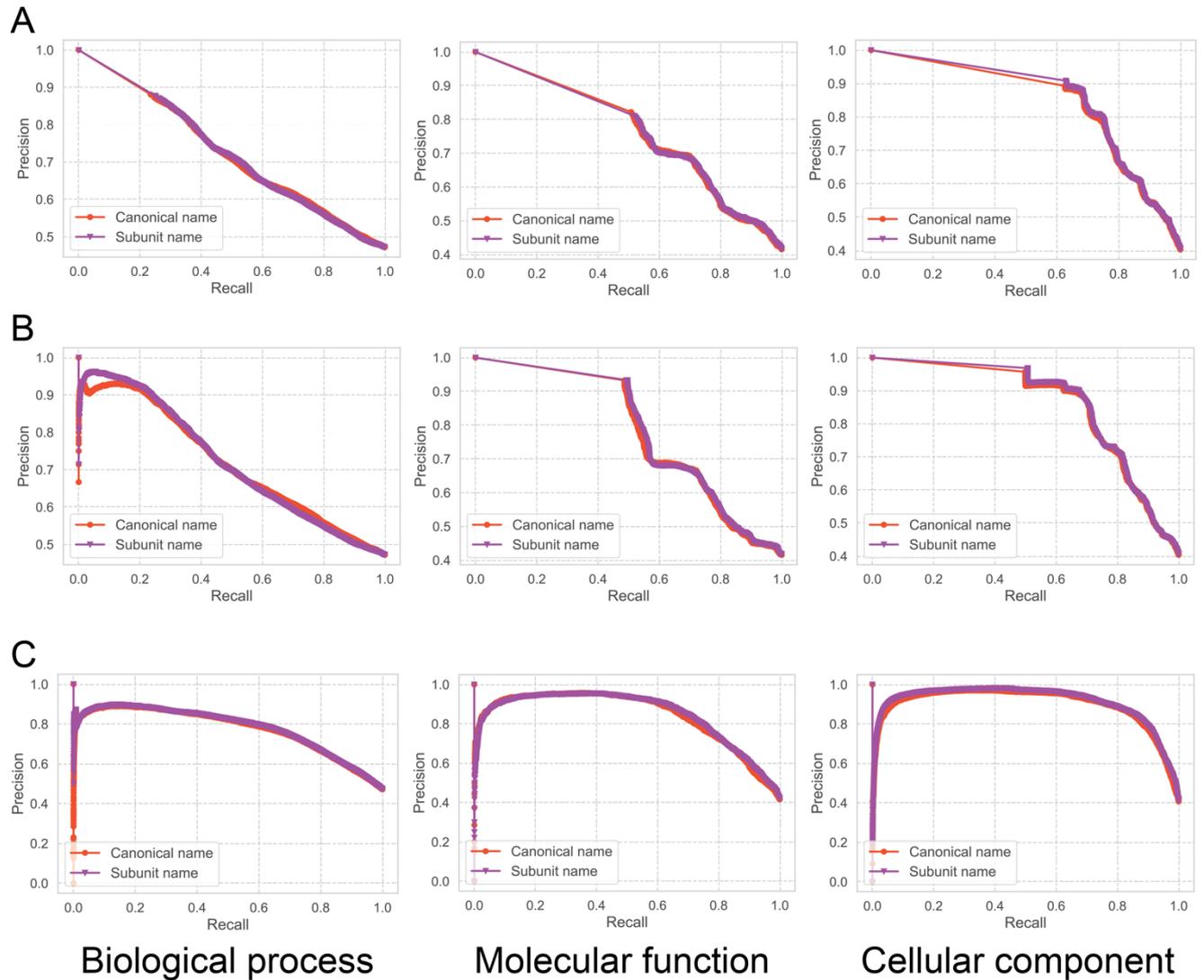
¹Acc: Accuracy

350

351

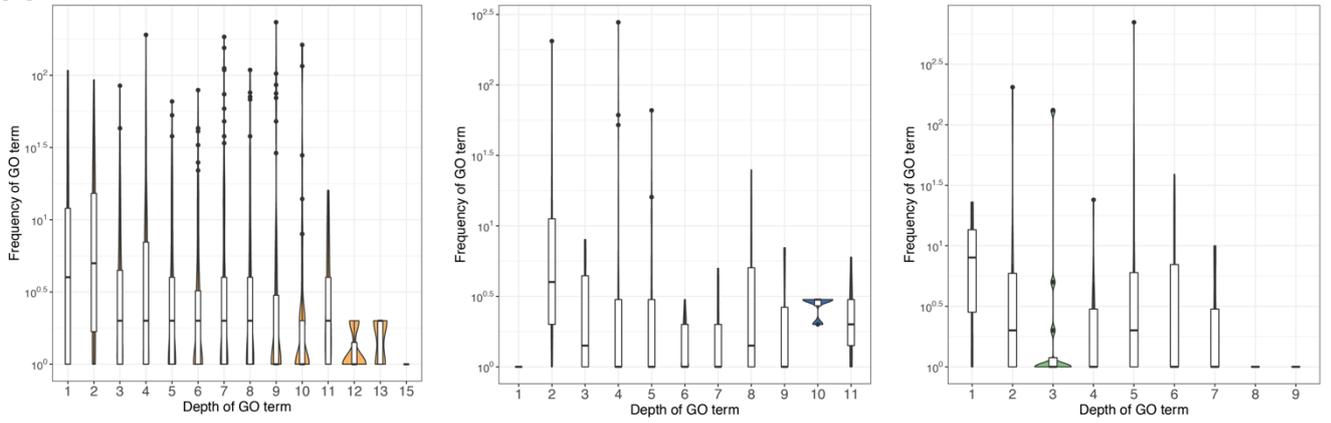


852 **Biological process** **Molecular function** **Cellular component**
 853 **Figure S1.** ROC curves and the average AUC values of (A) NB_Gauss, (B) NB_Bernoulli and (C) LR
 854 models on biological process, molecular function and cellular component categories via the adapted
 855 protein complex leave-one-out cross-validation. These models were trained using the CORUM ground-
 856 truth PC-GO associations, where the protein complexes were represented using both canonical and
 857 subunit naming schemes.

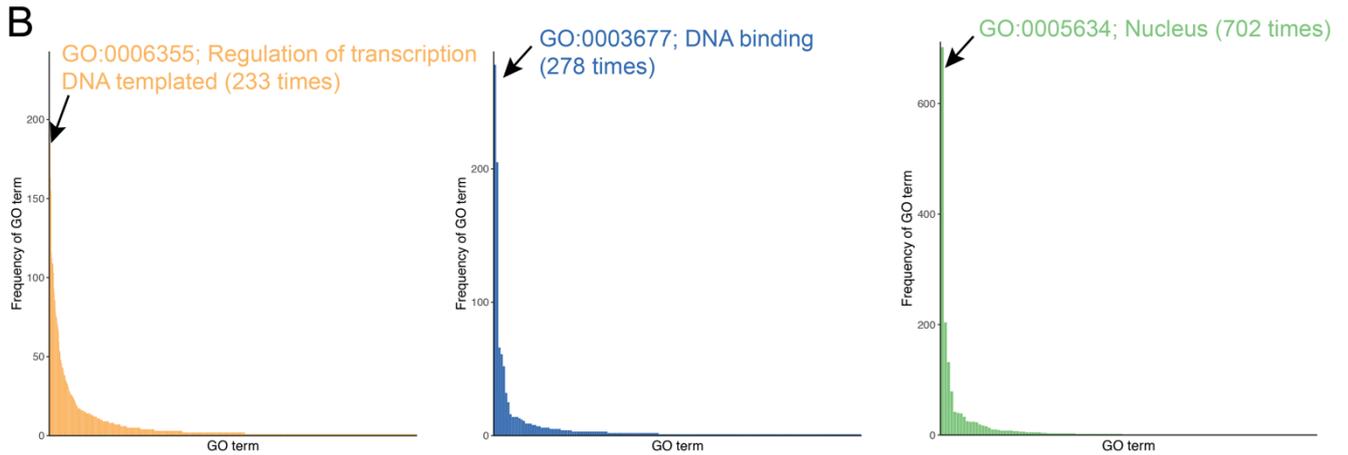


858
859 **Figure S2.** Precision-recall curves of (A) NB_Gauss, (B) NB_Bernoulli and (C) LR models on biological
860 process, molecular function and cellular component categories via the adapted protein complex leave-
861 one-out cross-validation, where the protein complexes were represented using both canonical and
862 subunit naming schemes.
863

A



B



Biological process

Molecular function

Cellular component

864
 865 **Figure S3.** Statistical analyses of the annotated GO terms in the CORUM database, including (A) the
 866 distributions of depth of annotated biological process, molecular function and cellular component terms
 867 in the GO DAG structures and (B) the frequencies of biological process, molecular function and cellular
 868 component terms that were assigned to the CORUM protein complexes. The top over-annotated
 869 biological process, molecular function and cellular component terms are respectively indicated.