

A peer-reviewed version of this preprint was published in PeerJ on 2 August 2018.

[View the peer-reviewed version](https://doi.org/10.7717/peerj.5298) (peerj.com/articles/5298), which is the preferred citable publication unless you specifically need to cite this preprint.

Doğan T. 2018. HPO2GO: prediction of human phenotype ontology term associations for proteins using cross ontology annotation co-occurrences. PeerJ 6:e5298 <https://doi.org/10.7717/peerj.5298>

1 **HPO2GO: Prediction of Human Phenotype Ontology Term**
2 **Associations Using Cross Ontology Annotation Co-occurrences**

3

4 Tunca Doğan^{1,2,3,*}

5

6 ¹Cancer Systems Biology Laboratory (CanSyL), Graduate School of Informatics, METU,
7 Ankara, 06800, Turkey

8 ²Department of Health Informatics, Graduate School of Informatics, METU, Ankara, 06800,
9 Turkey

10 ³European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI),
11 Hinxton, Cambridge, CB10 1SD, UK

12 * corresponding author email address: tdogan@metu.edu.tr

13

14 **ABSTRACT**

15 Analysing the relationships between biomolecules and the genetic diseases is a highly active
16 area of research, where the aim is to identify the genes and their products that cause a
17 particular disease due to functional changes originated from mutations. Biological ontologies
18 are frequently employed in these studies, which provided researchers with extensive
19 opportunities for knowledge discovery through computational data analysis.

20 In this study, a novel approach is proposed for the identification of relationships between
21 biomedical entities by automatically mapping phenotypic abnormality defining HPO terms
22 with biomolecular function defining GO terms, where each association indicates the
23 occurrence of the abnormality due to the loss of the biomolecular function expressed by the
24 corresponding GO term. The proposed HPO2GO mappings were extracted by calculating the
25 frequency of the co-annotations of the terms on the same genes/proteins, using already
26 existing curated HPO and GO annotation sets. This was followed by the filtering of the
27 unreliable mappings that could be observed due to chance, by statistical resampling of the
28 co-occurrence similarity distributions. Furthermore, the biological relevance of the finalized
29 mappings were discussed over selected cases, using the literature.

30 The resulting HPO2GO mappings can be employed in different settings to predict and to
31 analyse novel gene/protein - ontology term - disease relations. As an application of the
32 proposed approach, HPO term – protein associations (i.e., HPO2protein) are predicted. In
33 order to test the predictive performance of the method on a quantitative basis, and to compare
34 it with the state-of-the-art, CAFA2 challenge HPO prediction target protein set was
35 employed. The results of the benchmark indicated the potential of the proposed approach, as
36 HPO2GO performance was among the best ($F_{max} = 0.35$). The automated cross ontology
37 mapping approach developed in this work can easily be extended to other ontologies as well,
38 to identify unexplored relation patterns at the systemic level. The datasets, results and the
39 source code of HPO2GO are available for download at: <https://github.com/cansyl/HPO2GO>.

40 1. INTRODUCTION AND BACKGROUND

41 Systematic definition of biomedical entities (e.g., diseases, abnormalities, symptoms, traits,
42 gene and protein attributes, activities, functions and etc.) is crucial for computational studies
43 in biomedicine. Ontological systems, composed of standardized controlled vocabularies, are
44 employed for this purpose. Human Phenotype Ontology (HPO) system annotates disease
45 records (i.e., terms and definitions about diseases, recorded in relevant databases) with a
46 standardized phenotypic vocabulary (Robinson *et al.*, 2008; Köhler *et al.*, 2017). The source
47 of the disease information for HPO are Orphanet (Rath *et al.*, 2012), DECIPHER (Firth *et*
48 *al.*, 2009), and OMIM (Amberger *et al.*, 2014) databases. Each of the phenotype terms define
49 a specific type of abnormality encountered in human diseases (e.g., HP:0001631 - atrial
50 septal defect). The generation of HPO terms and their associations with diseases are carried
51 out with both manual curation efforts and automated procedures (e.g., text mining). The
52 curation job is usually done by experts by reviewing the relevant literature publications along
53 with the disease centric information at various biomedical data resources. The growing
54 library of HPO currently contains nearly 12,000 phenotype terms, providing more than
55 123,000 annotations to 7,000 different rare (mostly Mendelian) diseases and the newly added
56 132,000 annotations to 3,145 common diseases (Groza *et al.*, 2015). A long-term goal of the
57 HPO project is that the system to be adopted for clinical diagnostics, which will both provide
58 a standardized approach to medical diagnostics and present structured machine readable
59 biomedical data for the development of novel computational methods using data mining
60 techniques. Apart from phenotype-disease associations, which is the main aim of the HPO
61 project, HPO also provides phenotype-gene associations by using the known rare disease -
62 gene relations (i.e., the information which is in the form of: "certain mutation(s) in *Gene X*
63 causes the hereditary *Disease Y*"), using the abovementioned disease centric resources. The
64 associations between HPO terms and biomolecules, together with the downstream analysis
65 of these associations, help in disease gene identification and prioritization (Köhler *et al.*,
66 2009). With the mapping of phenotypes to human genes, HPO currently (January 2018)
67 provides 122,166 annotations between 3,698 human genes and 6,729 HPO terms.

68 The Gene Ontology (GO) is an ontological system to define gene/protein attributes with an
69 extensive controlled vocabulary (GO Consortium, 2014). Each GO term defines a unique
70 aspect of biomolecular attributes. Similar to other ontological systems GO has a directed

71 acyclic graph (DAG) structure, where terms are related to each other mostly with “is_a” or
72 “part_of” relationships. GO is composed of three categories (i.e., aspects) in terms of the type
73 of the defined gene product / protein attribute such as: (i) molecular function – MF (i.e., the
74 basic function of the protein at the molecular level; e.g., GO:0016887 - ATPase activity), (ii)
75 biological process - BP (i.e., the high level process, in which the protein plays a role; e.g.,
76 GO:0005975 - carbohydrate metabolic process), and (iii) cellular component – CC (i.e.,
77 subcellular location, where the protein carries out its intended activity; e.g., GO:0016020 -
78 membrane). Similar to the other ontological systems, the basic way of annotating a gene or
79 protein with a GO term is the manual curation by reviewing the relevant literature. GO also
80 employs the concept of “evidence codes”, where all annotations are labelled with descriptions
81 indicating the quality of the source information used for the annotation (e.g., ECO:0000006
82 - experimental evidence, ECO:0000501 – IEA: evidence used in automatic assertion).
83 UniProt-GOA (Gene Ontology Annotation) database (Huntley *et al.*, 2015) houses an
84 extensive collection of GO annotations for UniProt protein sequence and annotation
85 knowledgebase records. In the UniProtKB/Swiss-Prot database (i.e., housing manually
86 reviewed protein entries with highly reliable annotation) version 2018_02, there are a total
87 of 2,850,015 GO term annotations for 529,941 protein records; whereas in
88 UniProtKB/TrEMBL database (i.e., housing mostly electronically translated uncharacterized
89 protein entries) version 2018_02, there are a total of 189,560,296 GO term annotations for
90 67,760,658 protein records. Most of the annotations for the UniProtKB/TrEMBL database
91 entries are produced by automated predictions (UniProt Consortium, 2017).

92 Due to the high volume of experimental research that (i) discover new associations between
93 biomolecules and ontological terms, and (ii) produce completely new and uncharacterized
94 gene/protein sequences; curation efforts are having hard time in keeping up with annotation
95 process. To aid manual curation efforts, automated computational methods come into play.
96 These computational methods exploit the approaches and techniques widely used in the fields
97 of data mining, machine learning and statistics, to produce probabilistic associations between
98 biomedical entities. Critical Assessment of Functional Annotation (CAFA) challenge
99 (Radivojac *et al.*, 2013; Jiang *et al.*, 2016) aims to evaluate the automated methods that
100 produce GO and HPO term association predictions for protein entries, on a standard time-
101 held benchmarking dataset. Now after its third instalment, CAFA organization have already

102 brought together a research community, dedicated to elevate the capabilities of automated
103 function prediction approaches closer to the level of expert review.

104 Protein function prediction using GO terms is a highly active area of research, where various
105 types of approaches utilizing: amino acid sequence similarities (Hawkins *et al.*, 2009), 3D
106 structure analysis (Roy, Yang & Zhang, 2012), semantic similarities between the ontological
107 terms (Falda *et al.*, 2012), gene expression profiles (Lan *et al.*, 2013), protein-protein
108 interactions - PPIs (Wass, Barton & Sternberg, 2012), shared functional domains and their
109 arrangements (Fang & Gough, 2012; Finn *et al.*, 2016; Doğan *et al.*, 2016) and ensemble
110 approaches that exploit multiple feature types (Wass, Barton & Sternberg, 2012; Cozzetto *et*
111 *al.*, 2013; Lan *et al.*, 2013; Rifaioglu *et al.*, 2017); are employed to model the proteins and to
112 transfer the functional annotations from characterized proteins (i.e., the ones that have
113 reliable annotation), to the uncharacterized ones with highly similar features. Known GO
114 associations of genes and proteins are also used in different contexts in the literature. For
115 example, the method "MedSim" uses the semantic similarities between GO terms for the
116 prioritization of disease genes (Schlicker, Lengauer & Albrecht, 2010). The method "spgk"
117 uses a shortest-path graph kernel to compute functional similarities between gene products
118 using their GO annotations and the term relations on the GO DAG (Alvarez, Qi & Yan,
119 2011).

120 Apart from the machine-produced functional predictions for genes/proteins, automated
121 prediction of the associations between human genes/proteins and phenotype/disease defining
122 ontological terms is a non-trivial task, which can be utilized to identify large-scale novel
123 disease-gene-pathway/system relations. The identification of direct disease-gene relations is
124 a widely studied topic (Moreau & Tranchevent, 2012). A considerable amount of the existing
125 literature about disease-gene associations involve the calculation of semantic similarities
126 between gene products, based on the already existing ontological term annotations
127 (Washington *et al.*, 2009; Smedley *et al.*, 2013; Deng *et al.*, 2015; Rodríguez-García *et al.*,
128 2017). For example, the method "PhenomeNET" was employed to generate mappings
129 between the highly related terms across similar ontological systems (Rodríguez-García *et al.*,
130 2017) such as the HPO, Mammalian Phenotype Ontology – MP (Smith, Goldsmith & Eppig,
131 2005), Human Disease Ontology – DO (Kibbe *et al.*, 2014) and Orphanet Rare Disease
132 Ontology – ORDO (Vasant *et al.*, 2014); for discovering novel gene-disease associations.

133 However, semantic similarity based approach sometimes suffers from the low coverage of
134 especially the HPO annotations on the protein space. The authors of two recent studies have
135 investigated this issue (Kulmanov & Hoehndorf, 2017; Peng *et al.*, 2017). In this context,
136 increasing the coverage of HPO annotations by predicting gene/protein-HPO term
137 associations may help semantic similarity based association studies.

138 There are only a few examples of HPO term-protein association prediction methods in the
139 literature. In the "dcGO" method, the authors mapped ontological terms (including HPO) to
140 protein domains, which are the functional units, and transferred the ontology mapping to
141 proteins according to known domain annotations (Fang & Gough, 2012). The objective in
142 the "PHENOstruct" method is the prediction of gene-HPO term associations using
143 heterogeneous biological data consist of PPIs, GO annotations, literature relations, variants
144 and known HPO annotations, together with a structured SVM classifier (Kahanda *et al.*,
145 2015). One of the text mining based CAFA2 challenge participating methods "EVEX", was
146 employed for protein-HPO term association prediction. Originally, EVEX utilizes text
147 mining approaches for large-scale integration of heterogeneous biological data and event
148 extraction to generate a structured resource of relations, to be used in pathway curation (Van
149 Landeghem *et al.*, 2013). In the context of HPO term prediction, EVEX scans the literature
150 to detect proteins and phenotypic terms that co-occur on the same text corpus, and associates
151 them with each other based on certain criteria, similar to other text mining based approaches.
152 A network based HPO prediction method participated in CAFA2 was the "RANKS", in
153 which the authors developed a flexible algorithmic scheme for heterogeneous biological
154 network analysis, and used previously generated functional Interaction and functional human
155 gene networks for gene-HPO term association prediction (Valentini *et al.*, 2016). According
156 to the CAFA2 challenge results (Jiang *et al.*, 2016), the participating methods EVEX,
157 RANKS, PHENOstruct and dcGO were among the top performers. In a recent study, the
158 authors proposed two hierarchical ensemble methods: (i) the Hierarchical Top-Down, and
159 (ii) the True Path Rule, for gene-HPO term associations; in which the hierarchical graph
160 structure of HPO has been utilized together with the RANKS algorithm and the SVM
161 classifier (Notaro, *et al.*, 2017).

162 The text mining approach is highly effective for predicting gene-disease relations in disease
163 gene prioritization studies (Krallinger, Valencia & Hirschman, 2008). However, this

164 approach suffers from low coverage in some cases, due to knowledge limitation in the
165 literature. In other words, there is a bias towards detecting highly studied and already known
166 relations. If a certain abnormality and a gene/protein has not been studied together in the
167 same concept yet, it is often not possible to identify the relation. Network based methods are
168 proposed on top of either text-mining results, protein-protein interactions and/or pathway
169 data (Bromberg, 2013; Guney & Oliva, 2014; Guala & Sonnhammer, 2017) to detect indirect
170 relations, which greatly increase the coverage; nevertheless, they still moderately rely on the
171 previously reported relations. It is also important to note that, any predictive approach is
172 limited by the quality and the coverage of its source information. However, the predictive
173 output of different approaches often complement each other, contributing to fill different
174 portions of the missing information in the knowledge space. Due to this reason, developing
175 novel approaches to complement text mining based methods is crucial for automated
176 ontological association prediction. The observed low performance of even the best methods
177 in the HPO term prediction track of the CAFA2 challenge displayed the necessity of novel
178 approaches for the biomedical entity relation prediction.

179 In this study, a new approach is proposed to produce phenotypic abnormality HPO term
180 associations to both GO terms and human genes/proteins with the analysis of co-annotation
181 fractions between the HPO and GO term combinations. For this, HPO and GO terms that are
182 continually co-occurring on different proteins as annotations, are linked to each other (i.e.,
183 the system training step), entitled as the HPO2GO mappings. After that, proteins with a
184 linked GO term annotation receives the corresponding HPO term as the phenotypic term
185 prediction (i.e., the application step), entitled as the HPO2protein predictions. The idea here
186 is to associate a HPO term Y with a GO term X in the sense that: "if a protein loses its function
187 defined by the GO term X (or at least a reduction in the defined functionality) as a result of a
188 genetic mutation, the loss of function may cause the disease, which is defined by the
189 phenotype term Y ". This idea is based on the nature of annotating genes/proteins with HPO
190 terms; as for example, only the functionally problematic versions of these genes/proteins
191 (e.g., disease causing variants) are associated with the relevant genetic diseases and their
192 defining phenotypic abnormality terms. Mutations often lead to diseases by causing either a
193 loss of existing functionality or a gain of new functionality in the gene products. As a result,
194 if the HPO term Y and the GO term X are observed to be frequently co-occurring on different

195 proteins, then the lost function, which gave way to the corresponding disease may be the one
196 defined by the GO term X . This function usually corresponds to a large-scale biological
197 process. This approach exploits the significantly higher coverage of GO term annotations for
198 genes/proteins, compared to the HPO term annotations; to produce novel gene/protein - HPO
199 term associations.

200 In order to test the biological relevance of this approach, selected HPO2GO mappings were
201 manually examined. Additionally, the proposed methodology was employed to predict HPO
202 terms for the human protein target dataset provided in the CAFA2 challenge. Using the
203 benchmark set, the prediction performance was calculated and compared with the state-of-
204 the-art HPO prediction methods. Another set of HPO2GO mappings were generated for this
205 test, using the time-held training data provided in CAFA2. Finally, the up-to-date HPO2GO
206 mappings were employed to generate HPO term predictions to human protein entries in the
207 UniProtKB/Swiss-Prot database (i.e., HPO2protein predictions). The training and test
208 datasets, along with the source code of the proposed methodology and the analyses are
209 available for download at <https://github.com/cansyl/HPO2GO>.

210

211

212 2. METHODS

213 Dataset Construction

214 In order to generate the training sets, which were employed to generate the HPO2GO
215 mappings, first, gene to HPO term mappings file was downloaded from the HPO web-site
216 (January 2017 version of the file named: "ALL_SOURCES_ALL_FREQUENCIES_genes_
217 to_phenotype.txt"). This file contained 153,575 annotations between 3,526 human genes and
218 6,018 HPO terms. This file is shared in the HPO2GO repository with the filename:
219 "HPO_gene_to_phenotype_annotation_01_2017_ALL_SOURCES_ALL_FREQUENCIES
220 .txt". In HPO, "genes_to_phenotype" file only contains the asserted (i.e., specific)
221 annotations to genes; whereas "phenotype_to_genes" file contains all annotations propagated
222 through the root of the HPO DAG, according to the true path rule. As a result, parents of the
223 asserted terms are included as well. In this study, the asserted annotations are used in the
224 analysis (in terms of both GO and HPO), in order to make sure the training set includes only
225 the most reliable annotations.

226 Subsequently, all GO term annotations to the human proteins (with the experimental evidence
227 codes: EXP, IDA, IPI, IMP, IGI and IEP) in UniProtKB were downloaded from the UniProt-
228 GOA database 2017_01 version, using QuickGO browser (filename:
229 "GOA_UniProt_human_protein_annotation.tsv"). After eliminating the repeating (i.e.,
230 redundant) annotations, the finalized file contained 179,651 GO annotations between 18,577
231 unique human genes and 14,632 GO terms (filename of the finalized GO annotation file:
232 "GO_annot_human_proteins_UniProtGOA_01_2017.txt"). An additional column
233 containing the corresponding HGNC symbols (i.e., gene symbols) of the coding genes was
234 also included in the downloaded GO annotation file. This column was later used to combine
235 the GO annotations with the HPO annotations, since the HPO annotation file includes the
236 gene symbols.

237

238 Applied Methodology

239 The proposed methodology is divided into 2 steps: (i) training of the system (i.e., the
240 generation of the HPO2GO mappings), and (ii) the application step (i.e., the prediction of

241 HPO term-protein associations – HPO2protein, using the previously generated HPO2GO
242 mappings).

243 Figure 1 represents the whole HPO2GO mapping (i.e., training) procedure. For the training
244 of the system, first, the HPO and GO annotation datasets were prepared (Figure 1.1 and
245 Figure 1.2) and the initial HPO-GO mappings were generated (Figure 1.3) by identifying the
246 genes/proteins shared between individual HPO and GO terms (i.e., the cases where HPO and
247 GO terms are co-annotated to the same genes/proteins). This mapping generated 1,433,208
248 unique pairs between 6,005 HPO terms and 9,685 GO terms. At this point, it was observed
249 that some of GO and HPO terms were annotated to high number of proteins, and it was highly
250 probable to for them to co-occur on the same protein once or twice just by chance. In order
251 to eliminate the randomly occurred mapping cases, a filtering procedure was required to be
252 applied. For each HPO-GO term pair, a co-occurrence similarity measure, inspired from
253 semantic similarity based approaches, has been calculated. The co-occurrence similarity
254 formulation is given in Equation 1.

255

$$256 \quad S_{HPOi,GOj} = \frac{2 * N_{G\ HPOi\&\ GOj}}{N_{G\ HPOi} + N_{G\ GOj}} \quad (1)$$

257

258 Here, $S_{HPOi,GOj}$ is the co-occurrence similarity between the HPO term " $HPOi$ " and the GO
259 term " GOj ", $N_{G\ HPOi\&\ GOi}$ is the number of genes/proteins where these terms are annotated
260 together, $N_{G\ HPOi}$ is the total number of genes with the annotation " $HPOi$ ", and $N_{G\ GOi}$ is the
261 total number of genes with the annotation " GOi ".

262 The mapping process and the co-occurrence similarity calculation are shown in Figure 2 with
263 a toy example. Following the calculation of the co-occurrence similarities between all HPO-
264 GO pairs, a thresholding operation was applied in order to distinguish between relevant
265 mappings and the random ones. Two parameters were used for the thresholding operation:
266 (i) the co-occurrence similarities (S), and (ii) the number of genes with co-occurring
267 annotations (n). The aim behind employing a second parameter (i.e., n) was to eliminate the
268 potential random pairing cases, where the co-occurrence similarity is still high. These cases
269 are rare; however, it is still possible to observe a few of them especially when n is very small,
270 due to extremely high number of term combinations. In Figure 2, this situation is represented

271 on the toy example, here $S_{HPOD,GO4}$ is equal to $S_{HPOB,GO3}$; however the $HPOD-GO4$ mapping
272 is probably less reliable compared to $HPOB-GO3$ since $n_{HPOD,GO4}$ is equal to 1.

273 Statistical resampling was used to determine the optimal parameter values (to be used as
274 thresholds), that separate meaningful mappings from random ones. A permutation (i.e.,
275 randomization) test was constructed for this purpose. A randomized HPO-GO term mapping
276 table was generated (Figure 1.4) by first, shuffling the indices of the original "HPO vs. gene"
277 and "GO vs. gene" annotation tables; and second, calculating both the randomized co-
278 occurrence similarities (i.e., S_R) and the number of genes with co-occurring annotations (i.e.,
279 n_R) for each random HPO-GO mapping. For each arbitrarily selected S (i.e., $S > 0$, $S \geq 0.1$, S
280 ≥ 0.2 , ..., $S \geq 0.6$) and n (i.e., $n \geq 1$, $n \geq 2$, ..., $n \geq 5$) threshold value combination, the original
281 GO-HPO mappings with lower than the threshold S and n values were deleted and a co-
282 occurrence similarity distribution histogram was plotted using the remaining mappings (i.e.,
283 histograms plots in Figure 1 and in Figure 3). The same procedure was applied for the
284 randomized mapping set as well. Finally, Kolmogorov–Smirnov test -KS test- (Lilliefors,
285 1967; Hollander, Wolfe & Chicken, 2013) is employed to calculate a test statistic for
286 estimating whether the samples from the random and the original sets (at each S and n
287 selection) are from the same distribution or not. KS is a nonparametric test of 1-dimensional
288 probability distributions that can be used to compare two samples, considering the quantized
289 distance between the samples. The null hypothesis states that the two samples are drawn from
290 the same distribution. Here, the distribution (i.e., histogram) of the S values for the original
291 and the randomized mapping sets represent the two samples. The reason behind using
292 histograms instead of the actual S values was that, both high and low S values were presented
293 in both distributions; as a result, the significance test by checking sample distances approach
294 would not work. However, the frequencies of these high and low S values are different from
295 each other in the original and the random distributions. If the null hypothesis is accepted at a
296 selected threshold value pair (S and n), which means that the distributions are not statistically
297 different from each other, then it is concluded that the selected thresholds failed to eliminate
298 the random pairings in the original mapping (i.e., a higher threshold is required). The lowest
299 threshold values, where the samples from the two distributions became significantly different
300 from each other, were selected as the official thresholds. Excessive threshold values were not
301 considered in order not to eliminate too many GO-HPO mappings. After the determination

302 of the parameter values (i.e., S and n thresholds), the HPO2GO mappings were finalized,
303 which ended the training process.

304 HPO2protein prediction step was a simple procedure, where query proteins were annotated
305 with the HPO terms, by taking their already existing GO annotations into account. HPO2GO
306 mappings were employed for this purpose. There were a total of 3 application runs in this
307 study using: (i) CAFA2 targets as the query set (for the performance tests and for the
308 comparison with the state-of-the-art), (ii) CAFA3 targets as the query set (to officially
309 participate to the CAFA3 challenge, the results of which are yet to be announced), and (iii)
310 all human protein entries in the UniProtKB/Swiss-Prot database (to generate the
311 HPO2protein predictions).

312

313 **Performance Evaluation Metrics**

314 In this study, it was not possible to use a standard fold based cross-validation to measure the
315 performance and to determine the parameter values in the training procedure, since in most
316 cases, the number of genes/proteins that have a co-occurring HPO-GO term annotations were
317 so low. As a result, it was impossible to separate the samples into training and validation sets.
318 Instead, the optimal parameter values were determined by using statistical resampling.
319 However, a performance test was still required in order to assess the success of the proposed
320 approach. For this, CAFA2 challenge benchmark set was employed. Due to the fact that
321 CAFA2 challenge was long before the analysis done in this study, HPO2GO mappings were
322 re-generated using the training data provided in CAFA2. This was followed by the production
323 of the HPO-protein association predictions on the CAFA2 target gene set. This analysis both
324 served as a performance test with time-held data (one of the hardest and most informative
325 tests for predictive models) and a performance comparison with the state-of-the-art (i.e., other
326 HPO prediction methods participated in CAFA2). The most basic definitions of the
327 evaluation metrics used in this test; *recall*, *precision* and *Fmax*, are shown in Equation 2, 3
328 and 4.

329

330

$$Rc_{\tau i} = \frac{TP_{\tau i}}{TP_{\tau i} + FN_{\tau i}} \quad (2)$$

331

332

$$Pr_{\tau_i} = \frac{TP_{\tau_i}}{TP_{\tau_i} + FP_{\tau_i}} \quad (3)$$

333

334

$$F_{max} = \max_{i=1\dots N} \left\{ \frac{2 * Pr_{\tau_i} * Rc_{\tau_i}}{Pr_{\tau_i} + Rc_{\tau_i}} \right\} \quad (4)$$

335

336 In equations 2,3 and 4; TP_{τ_i} , FN_{τ_i} , FP_{τ_i} , Rc_{τ_i} and Pr_{τ_i} represent the number of true positives,
 337 the number of false negatives, the number of false positives, *recall* and *precision* values,
 338 respectively; at the i^{th} probabilistic score threshold. F_{max} correspond to the maximum of the
 339 *F-score* values (i.e., harmonic mean of *precision* and *recall*, shown inside the curly brackets
 340 in Equation 4) calculated for each arbitrarily selected probabilistic score threshold. Finally,
 341 $i=1\dots N$ represents there are N different arbitrarily selected probabilistic score thresholds.

342 In the proposed method, probabilistic scores for each HPO-protein association prediction is
 343 calculated using the term co-occurrence similarity scores in Equation 1. If the mapping
 344 between the terms HPO_i and GO_j received the co-occurrence similarity score S_{HPO_i,GO_j} , then
 345 all proteins that receive the HPO_i prediction due to the presence of GO_j annotation obtains
 346 the probabilistic prediction score: S_{HPO_i,GO_j} . The calculation of the score in Equation 1 is set
 347 to range between 0 and 1; as a result, it can directly be used as a probabilistic score. Apart
 348 from that, probabilistic score thresholds represent values, under which the predictions are
 349 discarded. This way, a different set of predictions are given for each arbitrarily selected
 350 probabilistic score thresholds, leading to different precision and recall values. It is important
 351 to note that, probabilistic score thresholds are different from the thresholds we used to filter
 352 out unreliable HPO2GO mappings during the training process. The probabilistic score
 353 thresholds are used here (i.e., after the production of HPO2protein predictions) to produce
 354 binary predictions from continuous prediction scores, to be able to calculate performances.
 355 More details regarding the CAFA2 evaluation metrics are given in Jiang *et al.*, 2016.

356

357 3. RESULTS

358 Statistical Analysis of the Mappings

359 The initial HPO to GO mappings were generated according to the procedure explained in the
360 Methods section (Figure 2). The initial mapping of the original set resulted in 1,433,208
361 mappings between 6,005 HPO terms and 9,685 GO terms. The same procedure for the
362 randomized set produced 1,543,917 mappings between 5,995 HPO terms and 9,685 GO
363 terms. The initial HPO-GO mappings for both the original and the randomized sets are
364 available for download in the repository of the study (respective filenames:
365 "HPO_GO_Raw_Original_Mapping.txt" and "HPO_GO_Random_Mapping.txt"). It was
366 expected that the mappings generated from the random set would have lower co-occurrence
367 similarity values on average compared to the original set mappings; in other words, they
368 would contain less number of mappings for a particular co-occurrence similarity value. Table
369 1 displays the comparison of the number of mappings for different co-occurrence similarity
370 values, between the original and the randomized sets. As observed from Table 1, when $S > 0$
371 there is no difference between the mappings; however as S is increased, the difference
372 between the mappings becomes clear. Also, when S is increased, the number of mapped HPO
373 and GO terms were decreased since many terms did not have any mappings that satisfied the
374 stringent S values. The parameter n was not taken into account while calculating the statistics
375 in Table 1 (i.e., $n \geq 1$ for all values in the table).

376 The histograms in Figure 3 display the co-occurrence similarity distributions (i.e., S) for
377 arbitrarily selected n values. As observed from the histograms, when the mappings with low
378 n values are eliminated, the distributions shift to the right (i.e., the mean of S increases),
379 which can be interpreted as the mappings became more reliable. However, excessive values
380 of n thresholds leave only a few mappings to work with, especially at $n=25$ and $n=75$ (please
381 see the number of mappings at the vertical axis of Figure 3.C and D). Histograms in Figure
382 3 also show that thresholding the mappings using only n (not using S at all) would not be
383 sufficient because there are mappings with very low S values even at very high n thresholds
384 (i.e., 25 and 75). This observation verified the decision to use both of the parameters for the
385 filtering operation. At this point, the statistical resampling (i.e., KS test) was applied since it

386 was not possible to determine the optimal n threshold by just manually checking the
387 histograms.

388 In order to find the minimum S and n values that significantly separate the original mapping
389 from the randomized mapping, 35 different distributions, all combinations of the selected n
390 (i.e., $n \geq 1, 2, \dots, 5$) and S (i.e., $S > 0, S \geq 0.1, \dots, 0.6$) values, were prepared and tested
391 individually against the co-occurrence distribution of the random mapping, generated with
392 the same S and n thresholds. This test resulted in 35 different p -value calculations and the
393 minimum parameter values that satisfied the statistical significance (i.e., rejection of the null
394 hypothesis, which states that the two samples are from the same distribution) were selected.
395 Table 2 displays the significance results of all KS tests. The cells with "NaN" indicate the
396 cases, where the test could not be completed due insufficient number of samples to calculate
397 the statistic. However, incomplete tests were not a problem since the aim here was observing
398 the minimum threshold values, where the distributions significantly diverge from each other
399 (NaNs are located far away from this point). In Table 2, the cell with the p -value written in
400 bold font (i.e., 0.0057) signifies the point, where the corresponding thresholds $n \geq 2$ and $S \geq$
401 0.1 yielded the required significance (p -value < 0.01); and thus, these values were selected
402 as the finalized thresholds. This means that, all of the mappings with $n < 2$ and $S < 0.1$ were
403 considered unreliable and eliminated from the initial HPO-GO mappings.

404 Figure 4 displays the total number of unique mappings (vertical axis) with co-occurrence
405 similarity values greater than the corresponding threshold value (horizontal axis), for the
406 original and the randomized distributions on the blue and red coloured curves, respectively.
407 Figure 4.A shows the plot for the combination with greater than or equal to one co-annotated
408 gene (i.e., $n \geq 1$), Figure 4.B displays the same value for $n \geq 2$, Figure 4.C and D for $n \geq 3$
409 and 4; respectively. The differences between Figure 3 and Figure 4 is that, (i) in Figure 4
410 cumulative number of mappings are given (i.e., all mappings left after thresholding with $S \geq$
411 0.1, 0.2, ...), whereas in Figure 3, the number of mappings that fall into each S bin is given;
412 and (ii) in Figure 4, plots are given for $n \geq 1, 2, 3$ and 4 since the aim was to display the
413 curves around the selected threshold n value; whereas in Figure 3, there are plots for $n \geq 1,$
414 5, 25 and 75 to visually indicate the distribution shifts especially at high n values (i.e., $n=25$
415 and $n=75$). Figure 4 was drawn as a visual representation of the likeness between the original
416 and the randomized distributions at different parameter selections. As observed from Figure

417 4, the distributions diverged from each other at $n \geq 2$, which also is consistent with the KS
418 test results. Considering the co-occurrence similarity parameter, $S \geq 0.1$ produced a clear
419 separation between the original and the randomized distributions as long as n is greater than
420 1. Following the HPO-GO mapping elimination according to the selected thresholds,
421 finalized HPO2GO mappings contained 45,805 associations between 3,693 HPO terms and
422 2,801 GO terms. HPO2GO mappings are available for download in the repository of the
423 study (filename: "HPO2GO_Finalized_Mapping.txt").

424 It was only possible to use a small portion of the input GO annotations for the generation of
425 the HPO2GO mappings because the number of HPO annotated genes were only 3,526;
426 whereas, the number of GO annotated human genes were 18,577. Since mappings can be
427 done over the genes/proteins with co-occurring GO and HPO annotations, only 3,526
428 genes/proteins were used in the process. The remaining 15,051 human genes with GO
429 annotations were only used in the application step (i.e., HPO2protein), to predict HPO term
430 associations.

431

432 **The Biological Relevance of the Selected HPO2GO Mappings – A Case Study**

433 Two different examples were selected and examined to discuss the biological relevance of
434 HPO2GO mappings. The first case is the mapping between the phenotypic abnormality HPO
435 term "absence of bactericidal oxidative respiratory burst in phagocytes" (HP:0002723) and
436 the GO term "respiratory burst after phagocytosis" (GO:0045730), which is in the BP
437 category. The exact definition of this GO term in the UniProt-GOA database is: "*A phase of*
438 *elevated metabolic activity, during which oxygen consumption increases; this leads to the*
439 *production, by an NADH dependent system, of hydrogen peroxide (H2O2), superoxide*
440 *anions and hydroxyl radicals*" (URL: <https://www.ebi.ac.uk/QuickGO/term/GO:0045730>).
441 These two terms are mapped to each other in HPO2GO with high confidence (i.e., $S = 0.89$
442 and $n = 4$). The symbols of the co-annotated genes were *CYBA*, *CYBB*, *NCF2* and *NCF1*. As
443 observed from the names of both terms and from the description of the GO term, the HPO
444 term defines an abnormal condition that corresponds to the absence of the biological process
445 portrayed by the mapped GO term. This is in accordance with the logic behind mapping HPO
446 terms with GO terms, which stated the occurrence of an abnormality (i.e., the HPO term) due

447 to the loss of the biomolecular function defined by the mapped GO term. There also is a GO
448 term named "respiratory burst after phagocytosis" (GO:0045728), which is related to the
449 mapped term (GO:0045730) on the GO DAG. This term (GO:0045728) defines a more
450 specific function that is the exact opposite of the mapped HPO term (HP:0002723),
451 semantically. There also is an evidence for the relation between HP:0002723 and
452 GO:0045728 in the OBO formatted term definitions of HPO (URL:
453 <http://purl.obolibrary.org/obo/hp.obo>). However, in HPO2GO, GO:0045728 could not be
454 mapped to HP:0002723 due to low coverage in the source GO annotation set. GO:0045728
455 was only annotated to one gene (symbol: *HCK*), which was not annotated to HP:0002723, as
456 a result, the mapping could not be generated. Nevertheless, the mapped GO term
457 (GO:0045730) still defined a sufficiently related function.

458 The second selected case was the mapping between the HPO term "cerebellar hemisphere
459 hypoplasia" (HP:0100307) and the MF category GO term "tRNA-intron endonuclease
460 activity" (GO:0000213). The exact definition of this specific GO term in the UniProt-GOA
461 database is: "*Catalysis of the endonucleolytic cleavage of pre-tRNA, producing 5'-hydroxyl
462 and 2',3'-cyclic phosphate termini, and specifically removing the intron*" (URL:
463 <https://www.ebi.ac.uk/QuickGO/term/GO:0000213>). These two terms were mapped to each
464 other in HPO2GO with high confidence (i.e., $S = 0.86$ and $n = 3$). The symbols of the co-
465 annotated genes were *TSEN2*, *TSEN34* and *TSEN54*. The HPO term HP:0100307 is
466 associated with the disease entry "Pontocerebellar Hypoplasia, Type 2C (PCH2C)"
467 (OMIM:612390) in the OMIM database. According to the disease definition, pontocerebellar
468 hypoplasia is a heterogeneous group of neurodegenerative disorders associated with
469 abnormally small cerebellum and brainstem, and the type 2C is characterized by a
470 progressive microcephaly from child birth (Barth, 1993). The occurrence of the disease is
471 associated with missense mutations in either *TSEN2*, *TSEN34* or *TSEN54* genes, which are
472 parts of the tRNA splicing endonuclease complex (Budde *et al.*, 2008). It was reported that,
473 due to the abovementioned mutations, there was a partial loss in the function of cleaving the
474 pre-tRNAs by the endonuclease complex (Budde *et al.*, 2008). This is another clear example
475 for a HPO term defining an abnormal condition, that is caused by the perturbation in the
476 function defined by the mapped GO term.

477

478 Performance Comparison with the State-of-the-art

479 The test for the comparison with the state-of-the-art had two objectives: (i) measuring the
480 performance of the method on a time-held dataset to observe the relevance of the proposed
481 approach, and (ii) investigating how the proposed method competes with the best performing
482 methods in the literature. For this, we have re-generated the HPO2GO mappings using the
483 CAFA2 training set, which contained 133,175 annotations between 5,586 HPO terms and
484 4,418 proteins, from October 2015. Whereas, CAFA2 evaluation set (i.e., benchmarking set)
485 contained 37,090 annotations between 2,838 HPO terms and 440 proteins. The reason behind
486 the presence of low number of annotations (and proteins) in the evaluation set was that, only
487 the HPO annotations produced between the time of the challenge participation deadline and
488 the end of the annotation collection period (a total duration of nearly 8 months) were used to
489 generate the time-held evaluation set. All of the datasets, the source code and the
490 supplementary files used in the CAFA2 challenge, and thus in this benchmarking experiment,
491 is available through the CAFA project repositories (URLs:
492 <https://github.com/yuxjiang/CAFA2> and <https://ndownloader.figshare.com/files/3658395>).

493 HPO2GO mappings generated using the CAFA2 training set contained 27,424 mappings
494 between 2,640 HPO terms and 2,488 GO terms. Considering the whole CAFA2 human target
495 protein set, this mapping produced 1,922,333 HPO predictions for 16,256 proteins and 2,640
496 HPO terms. The calculated performance of this prediction set was low ($F_{max} = 0.30$), mainly
497 due to high number of false positive (FP) hits. However, it is also probable that many of these
498 false positives were actually non-documented HPO annotations of the corresponding protein,
499 as the benchmark annotation set is incomplete. Increasing the thresholds with the aim of
500 reducing the number of false positives resulted in a matching increase in the number of false
501 negatives (FN), with a similar F_{max} value. With the aim of enriching the mappings (to be
502 able to reduce FPs without a significant increase in FNs), HPO annotations of genes from
503 January 2014 (i.e., the CAFA2 training set) were propagated to the root of HPO DAG
504 according to the true path rule. The propagated training set contained 379,513 annotations
505 between 4,418 human proteins and 6,576 HPO terms; as opposed to 133,175 annotations
506 between 4,418 human proteins and 5,586 HPO terms in the asserted CAFA2 set. As observed
507 from the dataset statistics, propagating the annotations have only added about one thousand
508 new terms to the set; however, the number of annotations were significantly increased.

509 Repeating the CAFA2 benchmark analysis using propagated HPO annotations and the same
510 GO annotations set resulted in the same performance ($F_{max} = 0.30$). Next, automated GO
511 annotations (i.e., evidence code: IEA) have been included in the source GO annotation set,
512 which increased the number of unique GO annotations from 128,947 to 214,235 (a 66%
513 increase). Using the propagated HPO annotations together with enlarged GO annotation set,
514 the new HPO-GO mappings, namely "HPOprop2GOall", were generated. The finalized
515 HPOprop2GOall contained 198,928 mappings between 4,780 HPO terms and 5,196 GO
516 terms; as opposed to 27,424 mappings between 2,640 HPO terms and 2,488 GO terms in the
517 original mappings. The drastic difference between the numbers have indicated the
518 enrichment provided by annotation propagation and GO set enlargement. Subsequently,
519 HPOprop2GOall mappings were used to predict HPO associations for all CAFA2 targets,
520 producing 13,022,574 predictions (as opposed to 1,922,333 predictions with the asserted set).
521 Considering only the CAFA2 benchmark proteins, the predictions generated by using the
522 optimized parameters (i.e., $n \geq 170$ and $S \geq 0.11$) resulted in 34,486 HPO predictions for 221
523 benchmark proteins and 235 HPO terms, with a performance of $F_{max} = 0.35$ (no-knowledge
524 benchmark sequences in the full evaluation mode), which is among the top performances
525 considering all of the models from 38 participating groups in the CAFA2 HPO prediction
526 track. The F_{max} performance of the top model in the challenge was 0.36 (Jiang *et al.*, 2016),
527 and the performance of the naïve baseline classifier was also the same. In Figure 5, each bar
528 displays the overall performance (F_{max}) of the CAFA2 participators, baseline classifiers and
529 HPO2GO. At this point in the study, additional HPO2GO mapping sets were generated using
530 different n and S threshold selections, and tested on the CAFA2 benchmark; however, these
531 mappings produced performances slightly inferior to the one generated using the optimal
532 thresholds (data not shown). HPO2GO CAFA2 benchmark test prediction results are
533 available in the repository of the study (filename:
534 "HPO_CAFA2_benchmark_predictions.txt").

535

536 **The Application of the Method to Generate Finalized HPO2protein Predictions**

537 Up-to-date HPO2GO mappings were employed to predict HPO terms for the human protein
538 entries in the UniProtKB/Swiss-Prot database (i.e., 20,258 protein records), and the resulting

539 prediction set was marked as the finalized HPO2protein predictions. This set contained
540 3,468,582 HPO predictions for 18,101 proteins and 3,693 HPO terms. HPO2protein
541 predictions are available in the repository of the study (filename:
542 "HPO2protein_Predictions.txt").

543 Finally, up-to-date HPO2GO model was run on the CAFA3 human protein targets, which
544 produced 3,453,130 predictions on 16,609 human proteins with 3,719 HPO terms. A more
545 stringent subset of this prediction set (i.e., predictions produced from mappings with $S \geq 0.2$)
546 has been officially submitted to the CAFA3 challenge. HPO2GO CAFA3 target predictions
547 are available in the repository of the study (filename:
548 "HPO_CAFA3_target_predictions.txt"). There was a small difference between the number
549 of query proteins in HPO2protein and the CAFA3 target sets (20,258 as opposed to 20,197,
550 respectively). At the time of writing this manuscript, the CAFA3 challenge results have not
551 been announced yet.

552

553 4. DISCUSSION

554 As a part of the main HPO project, a sub-set of the HPO terms had already been mapped to
555 the relevant terms from different ontology systems (e.g., anatomy, Gene Ontology process
556 or cell type) to yield semantic interoperability with these systems. However, this mapping
557 has been done by manually comparing the term definitions, only for a sub-set of GO terms;
558 as a result, the coverage of this mapping was quite limited. In our approach, we linked all
559 GO-HPO term combinations that satisfy the co-occurrence similarity tests. This way, the
560 non-documented relations are also identified. In this sense, it is expected that the HPO2GO
561 mappings will be valuable for the research community. It would also be interesting to
562 compare the HPO2GO mappings with the abovementioned manually curated associations;
563 however, it is not possible to access this data in the HPO repository anymore.

564 In this study, individual terms from both ontologies are mapped to each other considering the
565 co-annotated genes/proteins. However, the initial design of the experiment considered the
566 mapping of an HPO term to a trio of GO terms, one from each GO category (i.e., biological
567 process – BP, molecular function – MF and cellular component – CC). This way, the
568 corresponding phenotypic abnormality would be associated with a problem in a specific
569 molecular event (defined by the MF term), as a part of a defined large-scale process (BP
570 term), occurring at a particular sub-cellular location (CC term). This approach would have
571 been more biologically relevant compared to the current design; however, the initial design
572 failed due to the scarcity of both HPO annotations and GO annotations containing MF, BP
573 and CC term trios (data not shown). After that, a second option was considered, where HPO
574 terms were mapped to MF and BP term pairs to enrich the set of proteins with the required
575 GO annotations (i.e., MF and BP at the same time); nevertheless, the same problem was
576 encountered again. Reliable annotation sets with higher coverage, which may become
577 available in the future with more curation efforts, may solve this problem and make the
578 abovementioned mapping approach practical. However at present, even for the currently
579 applied one to one term mapping approach, the main challenge is the low coverage of the
580 predicted associations due to the small size of the source annotation sets. There can be a few
581 alternative solutions to this problem. First of all, the training sets with enriched GO
582 annotation (for the genes/proteins with HPO annotations) may be obtained by including the
583 annotations with evidence codes of reduced reliability (e.g., IEA – electronically generated).

584 Another option for enlarging the GO annotation set can be incorporating the genes (and their
585 respective annotations) from other organisms, that are orthologous to human genes. Scaling
586 up the coverage of HPO set can be provided by propagating the annotations to the parent
587 terms according to the hierarchical structure of HPO. Another option here would be taking a
588 more elaborate approach in the mapping procedure by taking the hierarchical term
589 relationships into account while generating the HPO2GO mappings (i.e., the parent and child
590 terms of the target HPO-GO term pair, that are co-annotated to different genes/proteins, will
591 also contribute to the calculation of the co-occurrence similarity of the target HPO-GO pair).

592 The official CAFA2 challenge results have indicated that, the methods based on sequence
593 similarities (e.g., the baseline classifier BLAST and a few models from the participating
594 groups) can achieve a good predictive performance considering the GO terms in the
595 molecular function (MF) category. This was expected since it is possible to detect most of
596 the signatures related to the molecular functions by analysing the amino acid sequence.
597 However, most of the sequence-similarity based methods failed in predicting the cellular
598 component (CC) GO term and HPO term associations. This can be explained for CC terms
599 as either by the cleavage of the signals from the sequence post-translationally or the
600 difficulties in detecting weak signals used for directing proteins to different compartments.
601 Considering the HPO prediction, the case may completely be different. As opposed to GO
602 terms, which define the attributes the proteins contain, HPO terms define phenotypic
603 abnormalities caused by the protein when it loses one (or more) of its functions, usually due
604 to certain mutations in the gene that codes the protein. Due to this reason, transferring a HPO
605 annotation from one protein to another based on sequence similarity does not have a
606 biological relevance, which explains the poor performance of the BLAST classifier.

607 An important observation regarding the CAFA tests done in this study is that, there was a
608 large difference between the number of HPO predictions for CAFA2 and CAFA3 targets,
609 using HPO2GO with default parameters (i.e., 1,922,333 in CAFA2 as opposed to 3,453,130
610 in CAFA3). There was also an increase in the number of predicted HPO terms (i.e., 2,640 in
611 CAFA2 as opposed to 3,719 in CAFA3), and there were no significant increase in the number
612 of targets. The increase in the number of predictions and the predicted HPO terms can be
613 attributed to the training set getting larger and more informative in time. The training set used
614 for CAFA2 contained 133,175 annotations; whereas, it was 153,575 for CAFA3. The

615 comparison of the predictive performances of HPO2GO trained by the CAFA2 and the
616 CAFA3 training sets may reveal more about the situation.

617

618 5. CONCLUSION

619 In this study, a simple and effective strategy, HPO2GO, is proposed to semantically map
620 phenotypic abnormality defining HPO terms with biomolecular function defining GO terms,
621 considering the cross-ontology annotation co-occurrences on different genes/proteins. This
622 approach can easily be translated into novel HPO term predictions for genes/proteins, as well
623 as into new HPO-disease or gene-disease associations. A literature based case study was
624 carried to discuss the biological relevance of the selected HPO2GO mappings. This work
625 also presents an application of the cross-ontology term mapping approach by generating
626 HPO-protein associations. HPO2GO was benchmarked on CAFA2 challenge protein targets
627 and it was revealed that the method was among the best performers of the HPO term
628 prediction track participators (i.e., the state-of-the-art methods). Also, the up-to-date trained
629 system was employed to predict HPO associations for all human proteins in the
630 UniProtKB/Swiss-Prot database (i.e., HPO2protein predictions). The methodology proposed
631 here was only meant to support the already established approaches (e.g., text mining), since
632 different techniques with different data sources and perspectives produce results that
633 complement distinct missing pieces of the knowledge space. It would also be interesting to
634 analyse the complementarity between the results of the proposed method and the results of
635 the conventional approaches participated in CAFA2 challenge; however, this was not
636 possible since the actual predictions of the participant groups are not publicly available.

637 As for the future work, it is first planned to map the HPO terms to GO term trios (i.e., MF,
638 BP and CC terms at the same time) using enriched annotation datasets, as explained at the
639 Discussion section. Another future task is the integration of HPO2GO mappings to our freely
640 available GO based automated protein function prediction tool/server UniGOPred (Rifaioglu
641 *et al.*, 2018); so that, query proteins that receive a GO term prediction will be automatically
642 associated with the HPO term(s) that are mapped to the corresponding GO term. It is expected
643 that this approach would produce large-scale HPO predictions for uncharacterized proteins
644 without any curated annotation, where the only available information is the amino acid
645 sequence. The knowledge extraction methodology proposed here can easily be combined
646 with various types of protein features employed in other predictive methods (e.g., variant
647 information, PPIs, gene expression profiles, etc.) to generate an ensemble HPO term
648 prediction tool that produces novel HPO-gene/protein-disease associations.

649 **6. REFERENCES**

- 650 Alvarez MA, Qi X, Yan C. 2011. A shortest-path graph kernel for estimating gene product
651 semantic similarity. *Journal of biomedical semantics*, 2(1), 3.
- 652 Amberger JS, Bocchini CA, Schiettecatte F, Scott AF, Hamosh A. 2014. OMIM. org: Online
653 Mendelian Inheritance in Man (OMIM®), an online catalog of human genes and genetic
654 disorders. *Nucleic acids research*, 43(D1), D789-D798.
- 655 Barth PG. 2014. Pontocerebellar hypoplasias: an overview of a group of inherited
656 neurodegenerative disorders with fetal onset. *Brain Dev.*, 15: 411-422.
- 657 Bromberg Y. 2013. Disease gene prioritization. *PLoS computational biology*, 9(4),
658 e1002902.
- 659 Budde BS, Namavar Y, Barth PG, Nürnberg G, Becker C, van Ruissen F, ... van der Knaap
660 MS. 2008. tRNA splicing endonuclease mutations cause pontocerebellar hypoplasia. *Nature*
661 *genetics*, 40(9), 1113.
- 662 Cozzetto D, Buchan DW, Bryson K, Jones DT. 2013. Protein function prediction by massive
663 integration of evolutionary analyses and multiple data sources. *BMC bioinformatics*, 14(3),
664 p.S1.
- 665 Deng Y, Gao L, Wang B, Guo X. 2015. HPOSim: an R package for phenotypic similarity
666 measure and enrichment analysis based on the human phenotype ontology. *PloS one*, 10(2),
667 e0115692.
- 668 Doğan T, MacDougall A, Saidi R, Poggioli D, Bateman A, O'Donovan C, Martin MJ. 2016.
669 UniProt-DAAC: domain architecture alignment and classification, a new method for
670 automatic functional annotation in UniProtKB. *Bioinformatics*, 32(15), 2264-2271.
- 671 Falda M, Toppo S, Pescarolo A, Lavezzo E, Di Camillo B, Facchinetti A, ... Fontana P. 2012.
672 Argot2: a large scale function prediction tool relying on semantic similarity of weighted Gene
673 Ontology terms. *BMC bioinformatics*, 13(4), S14.
- 674 Fang H, Gough J. 2012. DcGO: database of domain-centric ontologies on functions,
675 phenotypes, diseases and more. *Nucleic acids research*, 41(D1), D536-D544.
- 676 Finn RD, Attwood TK, Babbitt PC, Bateman A, Bork P, Bridge AJ, ... Gough J. 2016.

- 677 InterPro in 2017—beyond protein family and domain annotations. *Nucleic acids research*,
678 45(D1), D190-D199.
- 679 Firth HV, Richards SM, Bevan AP, Clayton S, Corpas M, Rajan D, ... Carter NP. 2009.
680 DECIPHER: database of chromosomal imbalance and phenotype in humans using ensembl
681 resources. *The American Journal of Human Genetics*, 84(4), 524-533.
- 682 Gene Ontology Consortium. 2014. Gene ontology consortium: going forward. *Nucleic acids*
683 *research*, 43(D1), D1049-D1056.
- 684 Groza T, Köhler S, Moldenhauer D, Vasilevsky N, Baynam G, Zemojtel T, ... Vasant D.
685 2015. The human phenotype ontology: semantic unification of common and rare disease. *The*
686 *American Journal of Human Genetics*, 97(1), 111-124.
- 687 Guala D, Sonnhammer EL. 2017. A large-scale benchmark of gene prioritization methods.
688 *Scientific reports*, 7, 46598.
- 689 Guney E, Oliva B. 2014. Analysis of the robustness of network-based disease-gene
690 prioritization methods reveals redundancy in the human interactome and functional diversity
691 of disease-genes. *PLoS one*, 9(4), e94686.
- 692 Hawkins T, Chitale M, Luban S, Kihara D. 2009. PFP: Automated prediction of gene
693 ontology functional annotations with confidence scores using protein sequence data.
694 *Proteins: Structure, Function, and Bioinformatics*, 74(3), 566-582.
- 695 Hollander M, Wolfe DA, Chicken E. 2013. *Nonparametric statistical methods*. John Wiley
696 & Sons.
- 697 Huntley RP, Sawford T, Mutowo-Meullenet P, Shypitsyna A, Bonilla C, Martin MJ, *et al.*
698 2015. The GOA database: gene ontology annotation updates for 2015. *Nucleic acids*
699 *research*, 43(D1):1057–63.
- 700 Jiang Y, Oron TR, Clark WT, Bankapur AR, D'Andrea D, Lepore R, ... Penfold-Brown D.
701 2016. An expanded evaluation of protein function prediction methods shows an improvement
702 in accuracy. *Genome biology*, 17(1), 184.
- 703 Kahanda I, Funk C, Verspoor K, Ben-Hur A. 2015. PHENOstruct: Prediction of human

- 704 phenotype ontology terms using heterogeneous data sources. *F1000Research*, 4.
- 705 Kibbe WA, Arze C, Felix V, Mitra E, Bolton E, Fu G, Mungall CJ, Binder JX, Malone J,
706 Vasant D, Parkinson H, Schriml LM . 2014. Disease ontology 2015 update: an expanded and
707 updated database of human diseases for linking biomedical knowledge through disease data.
708 *Nucleic acids research*, 43:1071–8.
- 709 Köhler S, Schulz MH, Krawitz P, Bauer S, Dölken S, Ott CE, ... Robinson PN. 2009. Clinical
710 diagnostics in human genetics with semantic similarity searches in ontologies. *The American*
711 *Journal of Human Genetics*, 85(4), 457-464.
- 712 Köhler S, Vasilevsky NA, Engelstad M, Foster E, McMurry J, Aymé S, ... Brudno M. 2016.
713 The human phenotype ontology in 2017. *Nucleic acids research*, 45(D1), D865-D876.
- 714 Krallinger M, Valencia A, Hirschman L. 2008. Linking genes to literature: text mining,
715 information extraction, and retrieval applications for biology. *Genome biology*, 9(2), S8.
- 716 Kulmanov M, Hoehndorf R. 2017. Evaluating the effect of annotation size on measures of
717 semantic similarity. *Journal of biomedical semantics*, 8(1), 7.
- 718 Lan L, Djuric N, Guo Y, Vucetic S. 2013. MS-k NN: protein function prediction by
719 integrating multiple data sources. *BMC bioinformatics*, 14(3), p.S8.
- 720 Lilliefors HW. 1967. On the Kolmogorov-Smirnov test for normality with mean and variance
721 unknown. *Journal of the American statistical Association*, 62(318), 399-402.
- 722 Moreau Y, Tranchevent LC. 2012. Computational tools for prioritizing candidate genes:
723 boosting disease gene discovery. *Nature reviews genetics*, 13(8), 523.
- 724 Notaro M, Schubach M, Robinson PN, Valentini G. 2017. Prediction of Human Phenotype
725 Ontology terms by means of hierarchical ensemble methods. *BMC bioinformatics*, 18(1),
726 449.
- 727 Peng J, Li Q, Shang X. 2017. Investigations on factors influencing HPO-based semantic
728 similarity calculation. *Journal of biomedical semantics*, 8(1), 34.
- 729 Radivojac P, Clark WT, Oron TR, Schnoes AM, Wittkop T, Sokolov A, ... Pandey G. 2013.
730 A large-scale evaluation of computational protein function prediction. *Nature methods*,

- 731 10(3), 221.
- 732 Rath A, Olry A, Dhombres F, Brandt MMC, Urbero B, Ayme S. 2012. Representation of rare
733 diseases in health information systems: the Orphanet approach to serve a wide range of end
734 users. *Hum. Mutat.*, 33:803–808.
- 735 Rifaioğlu AS, Doğan T, Saraç ÖS, Ersahin T, Saidi R, Atalay MV, Martin MJ, Cetin-Atalay
736 R. 2018. Large-scale automated function prediction of protein sequences and an experimental
737 case study validation on PTEN transcript variants. *Proteins: Structure, Function, and*
738 *Bioinformatics*, 86(2), 135-151.
- 739 Robinson PN, Köhler S, Bauer S, Seelow D, Horn D, Mundlos S. 2008. The Human
740 Phenotype Ontology: a tool for annotating and analyzing human hereditary disease. *The*
741 *American Journal of Human Genetics*, 83(5), 610-615.
- 742 Rodríguez-García MÁ, Gkoutos GV, Schofield PN, Hoehndorf R. 2017. Integrating
743 phenotype ontologies with PhenomeNET. *Journal of biomedical semantics*, 8(1), 58.
- 744 Roy A, Yang J, Zhang Y. 2012. COFACTOR: an accurate comparative algorithm for
745 structure-based protein function annotation. *Nucleic acids research*, 40(W1), W471-W477.
- 746 Schlicker A, Lengauer T, Albrecht M. 2010. Improving disease gene prioritization using the
747 semantic similarity of Gene Ontology terms. *Bioinformatics*, 26(18), i561-i567.
- 748 Smedley D, Oellrich A, Köhler S, Ruef B, Sanger Mouse Genetics Project, Westerfield M,
749 ... Mungall C. 2013. PhenoDigm: analyzing curated annotations to associate animal models
750 with human diseases. *Database*, bat025.
- 751 Smith CL, Goldsmith CAW, Eppig JT. 2005. The Mammalian Phenotype Ontology as a tool
752 for annotating, analyzing and comparing phenotypic information. *Genome biology*, 6(1), R7.
- 753 UniProt Consortium. 2017. UniProt: the universal protein knowledgebase. *Nucleic acids*
754 *research*, 45(D1), D158-D169.
- 755 Valentini G, Armano G, Frasca M, Lin J, Mesiti M, Re M. 2016. RANKS: a flexible tool for
756 node label ranking and classification in biological networks. *Bioinformatics*, 32(18), 2872-
757 2874.

- 758 Van Landeghem S, Björne J, Wei CH, Hakala K, Pyysalo S, Ananiadou S, ... Ginter F. 2013.
759 Large-scale event extraction from literature with multi-level gene normalization. *PloS one*,
760 8(4), e55814.
- 761 Vasant D, Chanas L, Malone J, Hanauer M, Olry A, Jupp S, ... Rath A. 2014. Ordo: An
762 ontology connecting rare disease, epidemiology and genetic data. *In Proceedings of ISMB*.
- 763 Washington NL, Haendel MA, Mungall CJ, Ashburner M, Westerfield M, Lewis SE. 2009.
764 Linking human diseases to animal models using ontology-based phenotype annotation. *PLoS*
765 *biology*, 7(11), e1000247.
- 766 Wass MN, Barton G, Sternberg MJ. 2012. CombFunc: predicting protein function using
767 heterogeneous data sources. *Nucleic acids research*, 40(W1), W466-W470.
- 768

769 **Table 1.** Statistics of the initial (i.e., raw) original and randomized HPO-GO mappings ($n \geq 1$).

S	# of mappings		# of mapped HPO terms		# of mapped GO terms	
	Original mapping	Random mapping	Original mapping	Random mapping	Original mapping	Random mapping
= 1	2 433	1 898	844	877	1 108	1 265
≥ 0.9	2 440	1 898	848	877	1 109	1 265
≥ 0.8	2 658	1 899	962	878	1 179	1 266
≥ 0.7	2 805	1 899	1 028	878	1 212	1 266
≥ 0.6	7 355	5 249	1 941	1 653	2 577	2 844
≥ 0.5	8 075	5 252	2 188	1 655	2 712	2 847
≥ 0.4	15 462	9 724	3 014	2 243	4 053	4 207
≥ 0.3	32 393	21 615	4 082	3 017	6 011	6 081
≥ 0.2	63 439	43 593	5 032	3 662	7 569	7 490
≥ 0.1	181 048	134 038	5 920	5 199	8 884	9 005
> 0.0	1 433 208	1 543 917	6 005	5 995	9 685	9 685

770

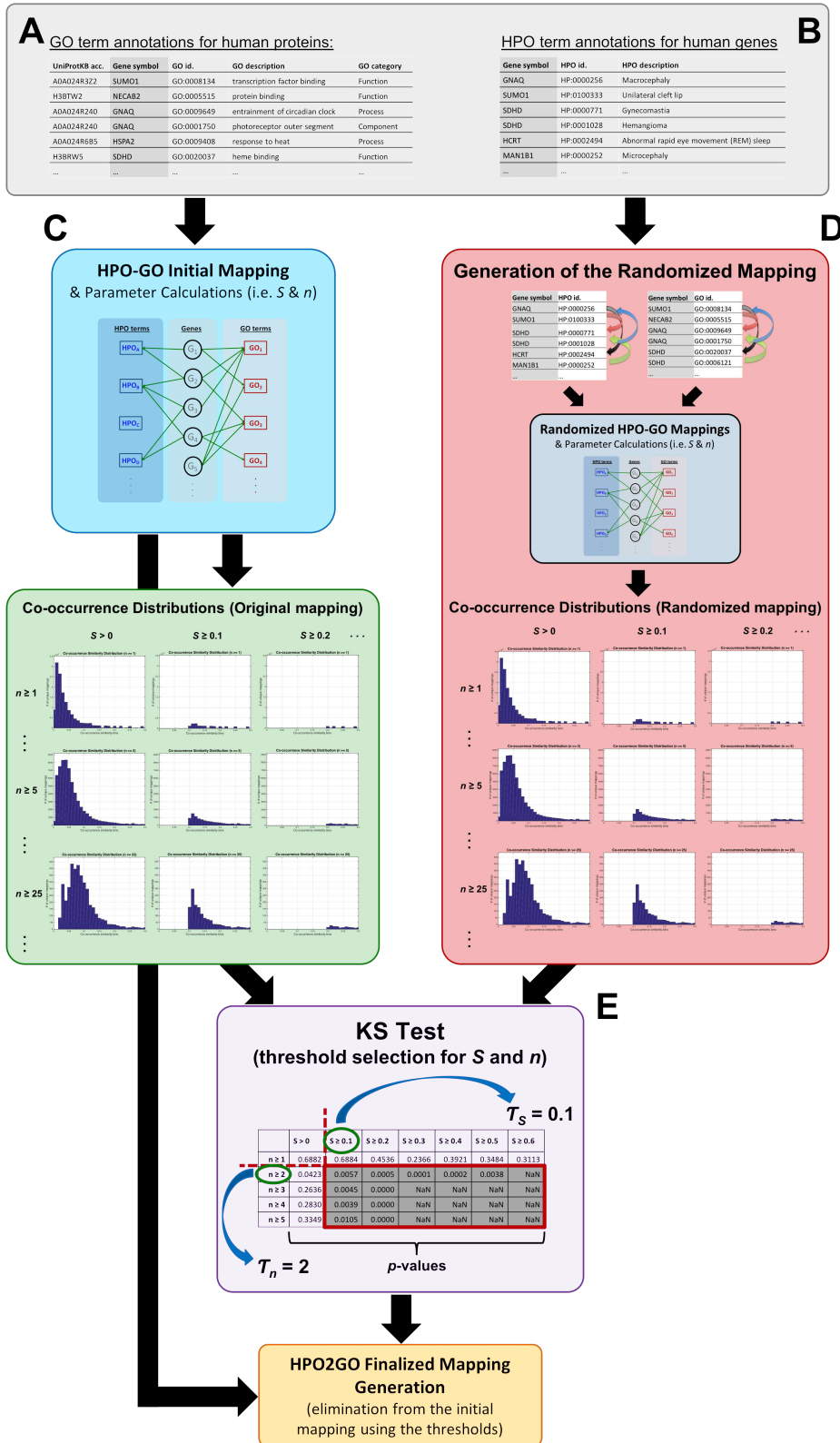
771

772 **Table 2.** KS test significance values for the comparison of original vs. randomized distributions at
 773 different co-occurrence similarity (S) and the number of co-annotated genes (n) thresholds.

KS test statistic		Co-occurrence similarity threshold						
		$S > 0$	$S \geq 0.1$	$S \geq 0.2$	$S \geq 0.3$	$S \geq 0.4$	$S \geq 0.5$	$S \geq 0.6$
# of co-annotated genes threshold	$n \geq 1$	0.6882	0.6884	0.4536	0.2366	0.3921	0.3484	0.3113
	$n \geq 2$	0.0423	0.0057	0.0005	0.0001	0.0002	0.0038	NaN
	$n \geq 3$	0.2636	0.0045	0.0000	NaN	NaN	NaN	NaN
	$n \geq 4$	0.2830	0.0039	0.0000	NaN	NaN	NaN	NaN
	$n \geq 5$	0.3349	0.0105	0.0000	NaN	NaN	NaN	NaN

774

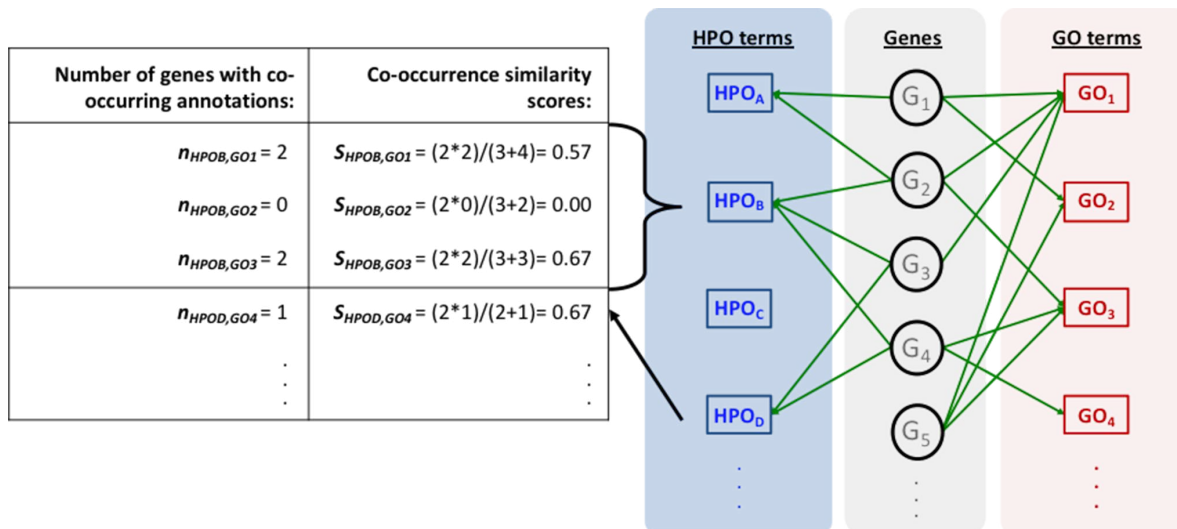
775



776

777 **Figure 1.** Schematic representation of the whole HPO2GO mapping (i.e., training) procedure.

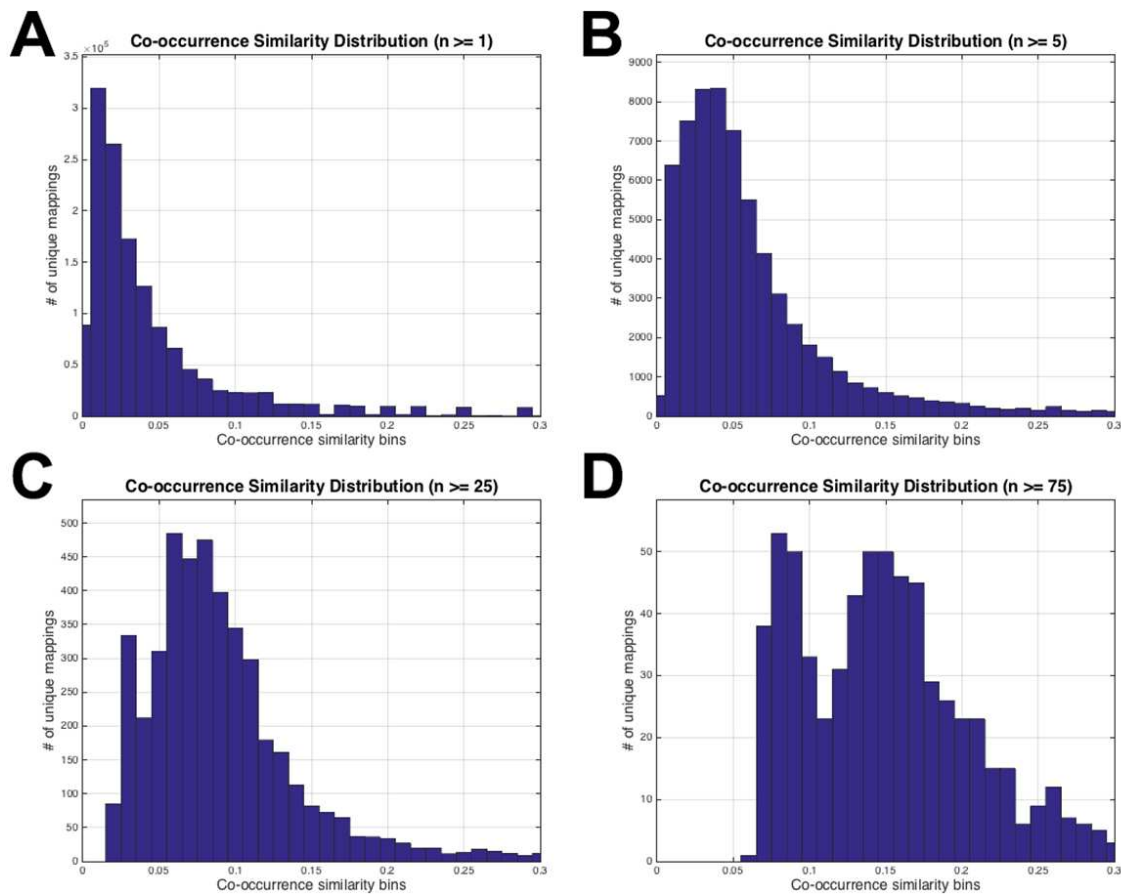
778



779

780 **Figure 2.** Representation of the initial HPO-GO mapping process together with the calculation of co-
 781 occurrence similarities (S) and the number of genes with co-occurring annotations (n), on a toy
 782 example.

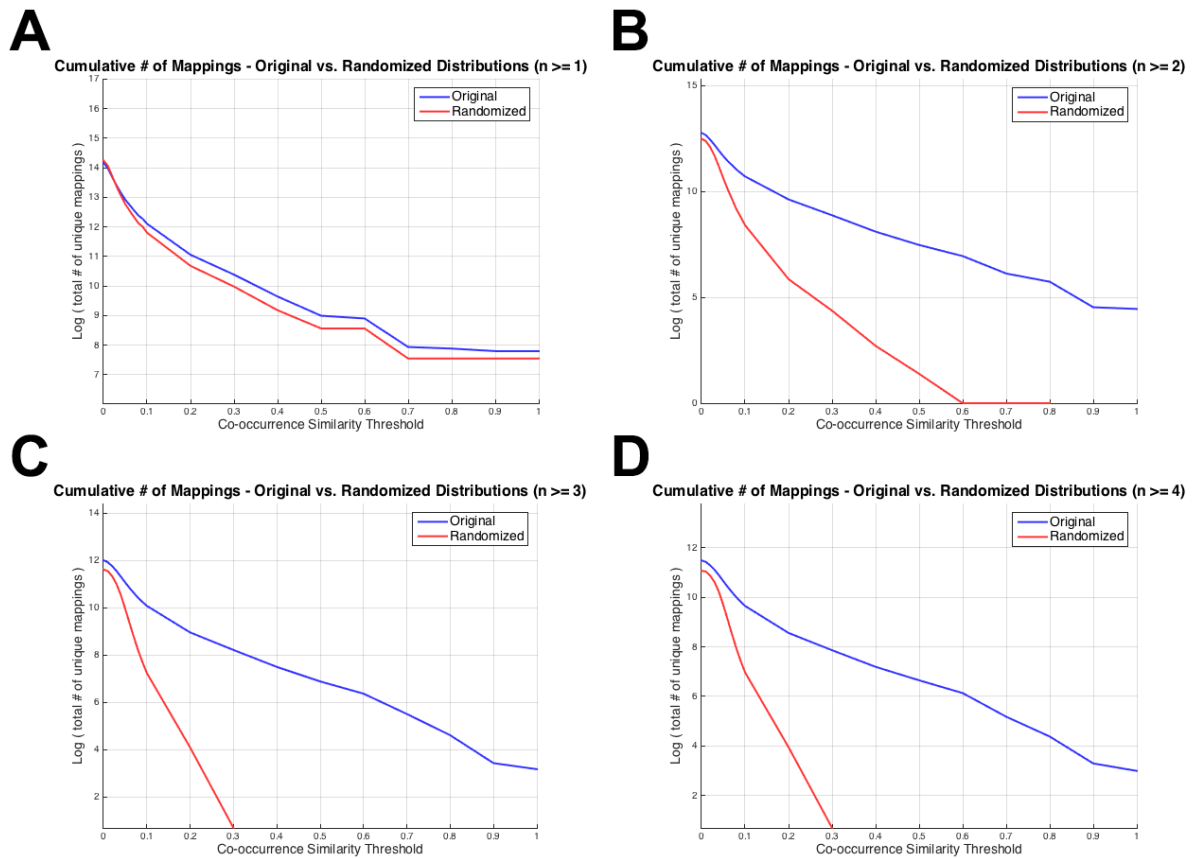
783



784

785 **Figure 3.** HPO-GO initial mappings co-occurrence similarity distributions. Each plot is drawn for a
786 different value of the number of co-annotated genes (i.e., n).

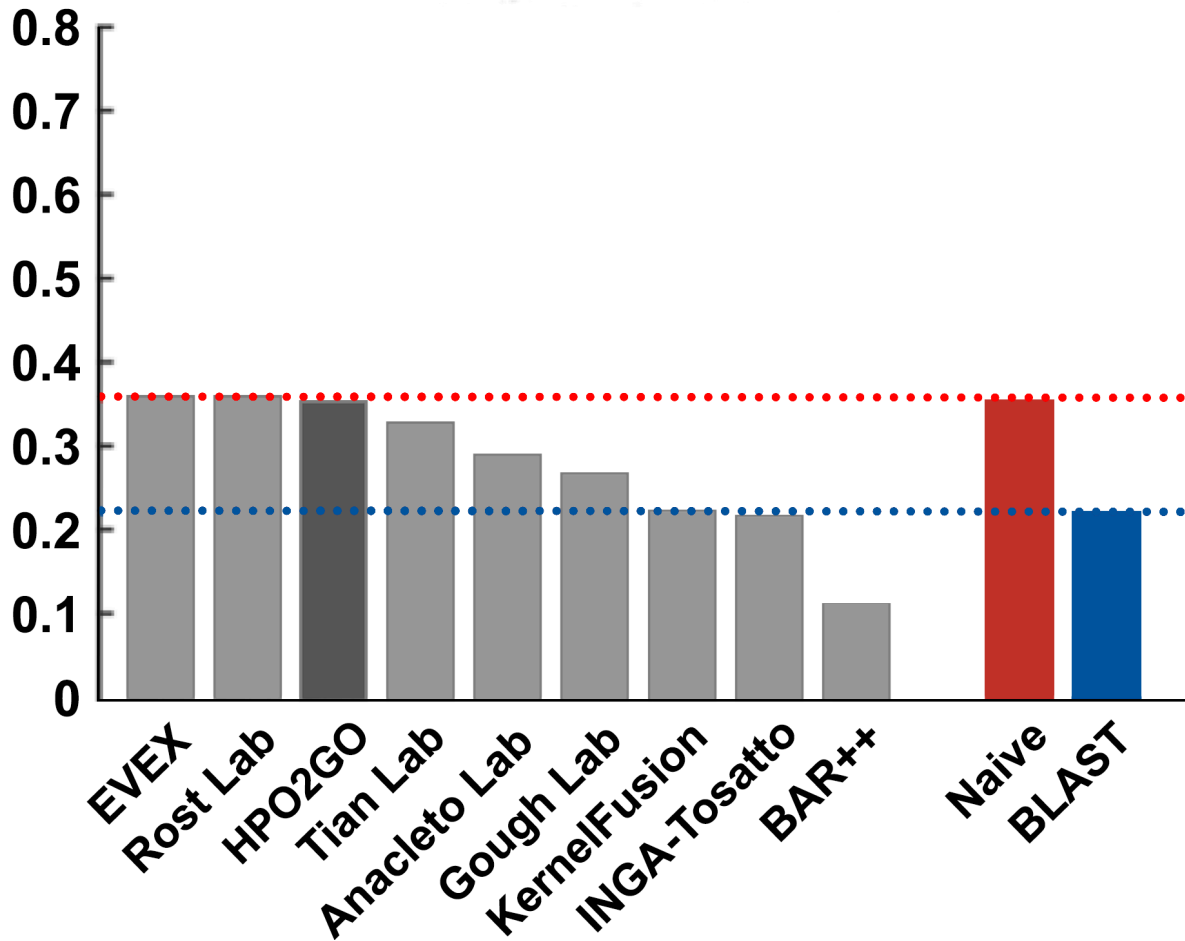
787



788

789 **Figure 4.** Cumulative plots displaying the number of HPO-GO mappings for the original (blue curve)
 790 and the randomized (red curve) distributions. Horizontal axis displays the arbitrarily selected co-
 791 occurrence similarity thresholds (i.e., τ_s), and the vertical axis represents the logarithm of the total
 792 number of mappings left after the application of the corresponding threshold. Each plot is drawn for
 793 a different value of the number of co-annotated genes (i.e., n). As the threshold (i.e., the minimum
 794 required co-occurrence similarity value to keep a mapping in the system) increase, more mappings
 795 are eliminated; thus, a monotonic decrease was observed for all plots.

796



797

798 **Figure 5.** F1-score performance results (F_{max}) of the top performing groups (grey bars), baseline
799 classifiers (red and blue bars) and HPO2GO (dark grey bar) in CAFA2 HPO prediction benchmark.
800 The lengths of the bars are directly proportional to the performance.

801