

1 Protein Function Prediction for newly sequenced organisms

2

3 Mateo Torres (mateo.torres@fgv.br)^{1,*},
4 Haixuan Yang (haixuan.yang@nuigalway.ie)^{2,*},
5 Alfonso E. Romero (aeromero@cs.rhul.ac.uk)^{3,*},
6 Alberto Paccanaro * (alberto.paccanaro@rhul.ac.uk)^{1,3}

7

8 1: Escola de Matemática Aplicada, Fundação Getúlio Vargas, Rio de Janeiro, Brazil

9 2: School of Mathematics and Statistical Sciences, National University of Ireland Galway,
10 Galway, Ireland

11 3: Department of Computer Science, Centre for Systems and Synthetic Biology, Royal
12 Holloway, University of London, Egham Hill, Egham, UK

13 *: These authors contributed equally to this work.

14

15 Abstract

16

17 Recent successes in protein function prediction have shown the superiority of approaches
18 that integrate multiple types of experimental evidence over methods that rely solely on
19 homology. However, newly sequenced organisms continue to represent a difficult
20 challenge, because only their protein sequences are available and they lack data derived
21 from large scale experiments.

22 We introduce S2F (Sequence to Function), a network propagation approach for the
23 functional annotation of newly sequenced organisms. Our main idea is to systematically
24 transfer functionally relevant data from model organisms to newly sequenced ones, thus
25 allowing us to use a label propagation approach. S2F introduces a novel label diffusion
26 algorithm that can account for the presence of overlapping communities of proteins with
27 related functions. Since most newly sequenced organisms are bacteria, we tested our
28 approach in the context of bacterial genomes. Our extensive evaluation shows a great

29 improvement over existing sequence-based methods, as well as four state-of-the-art
30 general-purpose protein function prediction methods.

31 Our work demonstrates that employing a diffusion process over networks of transferred
32 functional data is an effective way to improve predictions over simple homology. S2F is
33 applicable to any type of newly sequenced organism as well as to those for which
34 experimental evidence is available. A free, easy to run version of S2F is available at
35 <https://www.paccanarolab.org/s2f>.

36

37 **Introduction**

38

39 Less than 1% of the available protein sequences are currently annotated with reliable
40 information and the gap between unannotated and annotated sequences is widening at an
41 unprecedented rate¹ (Supplementary Note 1). Traditional experimental approaches to
42 determine protein function are usually expensive, time consuming, and provide low
43 throughput. While higher throughput approaches have recently been developed, they are
44 also proving to be insufficient to cope with the sheer number of new sequences produced
45 by next generation sequencing techniques². In this context, the computational annotation of
46 protein function has become a crucial step for a better understanding of the complex
47 mechanisms of living cells.

48

49 Newly sequenced organisms represent a particularly difficult challenge for automated
50 annotation methods because only their protein sequences are available and, in general, we
51 lack any other data derived from large scale functional experiments. In fact, protein function
52 prediction is somewhat easier for more studied organisms, including model organisms,
53 where multiple types of functional experimental evidence (e.g., gene expression,
54 proteomics data) are available that can be integrated with sequence information. The
55 Critical Assessment of Functional Annotation Challenge (CAFA)³ has indeed shown that
56 advanced methods that integrate multiple types of information for the prediction of Gene
57 Ontology (GO)⁴ terms significantly outperform methods that use only sequence information.

58

59 Network propagation approaches have been shown to be among the most successful
60 methods to predict protein function when some sort of experimental evidence is available⁵.

61 These methods combine and amplify existing knowledge about the function of some of the
62 proteins by propagating it through networks where nodes represent proteins, and edges
63 represent pairwise functional relationships between them that are derived from
64 experiments (e.g., physical interaction, co-occurrence in protein complexes, co-expression).
65 In other words, these methods expand an initial set of functional labels available for some
66 experimentally characterised proteins (seeds) to related neighbouring proteins, thus
67 exploiting the guilt-by-association principle, according to which highly connected nodes
68 should share similar functional properties. However, until now, these ideas could not be
69 applied to newly sequenced organisms, since in this case both the seeds and the networks
70 are unavailable.

71

72 This paper introduces S2F (sequence to function), a novel network propagation-based
73 method for the functional annotation of newly sequenced organisms. Our main idea is to
74 systematically transfer functionally relevant data that is available for model organisms to
75 newly sequenced organisms, thus allowing us to use network propagation to predict protein
76 function. S2F presents a novel network propagation algorithm that can account for the
77 presence of overlapping communities of proteins with related functions.

78

79 Since most newly sequenced organisms are bacteria, we have developed and tested our
80 solutions in the context of bacterial genomes. Bacteria is also the superkingdom with most
81 available sequenced proteins in UniProtKB (Supplementary Note 1), and the functional
82 characterisation of bacteria holds great potential in fields ranging from alternative energy
83 sources to understanding and treating disease. However, the ideas presented here are more
84 widely applicable to protein function prediction for any type of organism, and an earlier
85 version of our algorithm has successfully been applied to organisms from other kingdoms^{3,6}.

87 **Results**

88 The aim of S2F is to predict the function of each of the proteins in a newly sequenced
89 organism. Functional categories are defined according to the Gene Ontology (GO)⁴, where
90 terms are organised in a hierarchical structure with several domains and levels of specificity.
91 The prediction of protein function is a multi-class, multi-label classification problem: multi-
92 class, as there are over 40,000 possible GO terms that can be annotated to a protein; multi-
93 label, because each protein can be annotated with multiple GO terms. Importantly, the
94 hierarchical structure of the Gene Ontology must be taken into account for the prediction,
95 since whenever a protein is annotated with a GO term, it is also annotated with all its
96 ancestor terms up to the root of the ontology (this is known as the “true path rule”^{7,8}).
97 Therefore, an important requirement for the output of any protein function prediction
98 method is to be *consistent*: if a GO term is predicted with a certain probability, its parent
99 terms must be predicted with an equal or greater probability⁹.

100

101 S2F consists of four main components (see the pictorial representation in Fig. 1):

- 102 A. a method to *infer the initial seeds*, that combines the output of InterPro¹⁰ and
103 HMMER¹¹ to obtain a set of initial predictions that is consistent;
- 104 B. a method for *network transfer*, that relies on the concept of interolog^{12,13} to infer
105 several functional networks;
- 106 C. a method for *network combination*, that combines the different functional networks
107 into a single one;
- 108 D. a *label propagation* algorithm, that diffuses the seed information to obtain a
109 prediction.

110

111 In the following, we will describe each component in turn. We will assume that we wish to
112 predict the function for a newly sequenced organism (target organism) with n proteins, and
113 that the Gene Ontology contains t terms.

114

115 **A. S2F Seed Inference** InterPro¹⁰ constitutes an excellent starting point for predicting
116 protein function from sequence as it provides predictions from 14 different protein
117 signature databases. We consolidate its output into an $n \times t$ matrix of predictions R (see
118 Materials and Methods) which is consistent, and where each entry R_{ij} is the fraction of
119 InterPro models in which the (i, j) association is present.

120

121 Although InterPro predictions are extremely accurate, they are often limited in number and
122 involve only a few GO terms. In order to enrich the catalogue of GO terms that appear in our
123 initial seed set, HMMER¹¹ is run for every protein in the target organism against the
124 experimentally annotated sequences in UniProtKB/Swiss-Prot (Supplementary Note 2). This
125 results in the HMMER seed set, a binary matrix H of size $(n \times t)$, which is then up-
126 propagated according to the true path rule^{7,8}. A convex combination of H and R gives us the
127 consistent combined seed set $Y \in \mathbb{R}^{n \times t}$:

$$Y = \alpha R + (1 - \alpha)H$$

128 where $\alpha \in \mathbb{R}, 0 \leq \alpha \leq 1$ controls the relative contribution of InterPro and HMMER
129 predictions, and each entry of $Y, 0 \leq Y_{ij} \leq 1$.

130

131 **B. S2F Network Transfer** We build networks where nodes represent target organism
132 proteins, and edges represent pairwise functional relationships (interactions) between
133 them. Since experimental evidence of functional relationships between proteins is not

134 available for newly sequenced organisms, in order to create these networks we exploit the
135 fact that these relationships are often conserved across species^{14,15}. This allows us to
136 transfer existing evidence from well-studied organisms to newly sequenced ones.

137

138 Our starting point is the seminal work by Yu et al.¹³ who transferred different types of
139 functional networks with high precision using the concept of interolog-mapping first
140 proposed by Walhout et al.¹². The idea is that, given two proteins A and B in the target
141 organism, if there exists a pair of proteins A' and B' that are known to interact in another
142 organism (source organism), such that A is an orthologue of A', and B is an orthologue of B',
143 then we can infer an interaction between A and B.

144

145 Our transfer algorithm derives from the one proposed by Yu et al.¹³ (for details see
146 Materials and Methods and Supplementary Note 3). S2F uses STRING¹⁶ as the dataset of
147 different types of experimental interactions in source organisms. For each type of
148 interaction, S2F builds one transferred network, r , that can be represented as a matrix
149 $W^{(r)} \in \mathbb{R}^{n \times n}$, where each entry $W_{ij}^{(r)}$ represents the strength of the interaction between
150 proteins i and j in r . For a given target organism, S2F transfers five types of interaction,
151 namely “neighborhood”, “experiments”, “coexpression”, “textmining”, and “database”
152 using the experimental interactions available for any organism in STRING.

153

154 **C. S2F Network Combination** Having obtained a set of transferred networks for the target
155 organism, we now face the task of combining them into a single network for diffusing the
156 seeds. Our approach is to linearly combine the different networks through *learned*
157 coefficients. These coefficients provide us with interesting information about the relative
158 importance and role of each network in the prediction. While other systems learn this

159 combination (e.g., GeneMANIA¹⁷), the solution we propose here is applicable to our
160 problem, where no initial set of known labels is available.

161

162 We begin by using the InterPro predictions to build a network of functional similarities
163 $T \in \mathbb{R}^{n \times n}$, where the similarity between proteins i and j , T_{ij} , is defined as:

164

$$T_{ij} = \frac{|N_i \cap N_j|}{|N_i \cup N_j|}$$

165

166 where N_i and N_j are the sets of all GO terms above a threshold τ that are associated to
167 proteins i and j respectively in R , that is, $N_i = \{k | R_{ik} > \tau\}$, and $N_j = \{k | R_{jk} > \tau\}$.

168 Therefore, T_{ij} is the Jaccard similarity between sets of GO terms that are assigned by
169 InterPro to proteins i and j .

170

171 Given p networks $W^{(r)}$ with $r \in \{1, \dots, p\}$, we combine them into a single network
172 $W \in \mathbb{R}^{n \times n}$ using a weighted linear combination, where the vector of weights $\hat{c} \in \mathbb{R}^p$ is
173 learnt by minimising the square of the difference between T and the linear combination
174 (see Materials and Methods).

175

176 **D. S2F label propagation** Proteins rarely perform their functions in isolation, but rather they
177 act as part of functional groups. As mentioned earlier, network propagation methods for
178 protein function prediction exploit exactly this fact – groups of proteins that are highly
179 connected in functional networks form communities that share a similar function.
180 Importantly, when a protein has more than one function, it will belong to more than one of
181 such functional groups. We notice that such proteins, lying at the intersection of

182 communities are, in general, more functionally similar compared to their neighbours, since
183 they share more functional roles. Therefore, when a set of proteins has more than one
184 function, the propagation of information (or diffusion) between proteins within this set
185 should be higher than the diffusion between proteins in this set and proteins outside this
186 set. However, this does not happen with existing diffusion methods (for details see
187 Supplementary Note 6). Here we propose a novel label propagation method that explicitly
188 models overlapping communities and, in this way, corrects this problem.

189

190 We begin by defining the matrix $W^{S2F} \in \mathbb{R}^{n \times n}$, a transformation of the combined network
191 W whose entry W_{ij}^{S2F} is defined as:

192

$$W_{ij}^{S2F} = \frac{1}{2} \left(\frac{1}{d_i} + \frac{1}{d_j} \right) J_{ij} W_{ij}$$

193

194 where $d_i = \sum_{j=1}^n J_{ij} W_{ij}$ and J is a weighted Jaccard similarity matrix that models the
195 overlapping community effect (see Materials and Methods). We also define a diagonal
196 matrix D^{S2F} where the i -th diagonal element $D_{ii}^{S2F} = \sum_j W_{ij}^{S2F}$. Our algorithm produces a
197 prediction matrix $F \in \mathbb{R}^{n \times t}$ for all the n proteins of the organism and all the t GO terms by
198 computing the following:

199

$$F = (I + \lambda L)^{-1} Y$$

200

201 where Y is the matrix containing the initial labelling, I is the identity matrix, $L = D^{S2F} -$
202 W^{S2F} is the Laplacian of W^{S2F} , and $\lambda > 0$ is the regularisation parameter (see Materials
203 and Methods).

204

205 We show that this label propagation algorithm does not suffer from the problem described
206 above for overlapping communities (see Supplementary Note 6). Moreover, we prove that it
207 satisfies the necessary conditions to ensure that, for each pair of terms j and k such that j is
208 an ancestor of k (in these cases $Y_{ij} \geq Y_{ik}$ for every i), we have that $F_{ij} \geq F_{ik}$ for every i (the
209 proof can be found in Supplementary Note 7). As a consequence, since Y is consistent with
210 the Gene Ontology structure, F will also be consistent.

211 **Experimental Setup**

212 We present the evaluation of S2F on bacteria from UniProtKB. Following the evaluation
213 procedure used by most authors^{3,18} the performance of S2F in predicting protein function
214 was assessed both in a *per-gene* and in a *per-term* setting. In per-gene predictions, given a
215 gene, we assess the performance of S2F at predicting a set of functions associated to that
216 gene. Conversely, in per-term predictions, given a function, we assess the performance of
217 S2F at predicting a set of genes that perform that function.

218

219 The performance was assessed against a set of known experimental annotations. Therefore,
220 the bacteria used for testing were chosen so that they had at least a few experimentally
221 annotated genes (to be able to assess the performance in a per-gene setting) while
222 maintaining a reasonable diversity of annotated GO terms (to be able to assess the
223 performance in a per-term setting) in the GOA database¹⁹. The ten bacteria in Table 1
224 satisfied our set of criteria (the criteria are detailed in Materials and Methods).

225

226 This set of bacteria provides a good testbed for our experiments. The amount of
227 experimental annotations in these bacteria covers a wide spectrum, ranging from well-

228 studied bacteria (e.g., *E. coli*) to more obscure ones that are not even included in STRING
229 (e.g., *Brucella abortus*).

230

231 In our experiments, we tested the performance at predicting the functional annotation for
232 the whole genome for each of the ten bacteria, in turn. To avoid circular reasonings, when
233 testing each bacterium, we carefully removed any functional information for that bacterium
234 as well as for any phylogenetically close species. To do this, for each bacterium, we created
235 a list of excluded species in two steps. First, starting from that bacterium, we navigated the
236 NCBI taxonomy moving up two levels (i.e., to the parent of the parent node) and we
237 included in our list that node and all its descendants. Second, we added to the list all the
238 nodes in the NCBI taxonomy that had a similar name. Having created a list of excluded
239 species, we removed any information about these species from STRING, as well as about
240 their proteins from the GOA database. The detailed list of all organisms excluded when
241 testing each specific bacterium is provided in the Supplementary Data.

242

243 Predicted annotations were evaluated against the existing functional annotations (GOA files
244 in Supplementary Data) using the well-established metrics that have been used in the CAFA
245 challenge³: F_{\max} , S_{\min} , AUC-ROC, and AUC-PR metrics (for details, see Supplementary Note
246 12).

247

248 **Evaluation**

249 We compared the performance of S2F against InterPro, HMMER, Argot 2.5²⁰, DeepGOPlus²¹,
250 GOLabeler²² and NetGO²³. InterPro and HMMER are among the best and most widely used
251 sequence-based methods for predicting protein function for newly sequenced organisms.
252 The other four methods, although they were not explicitly conceived for this problem, could

253 nevertheless be employed here as they are able to predict protein function using sequence
254 information alone. Argot 2.5²⁰, and GOLabeler²² were among the top performer in the last
255 edition of the CAFA competition⁶; NetGO²³ and DeepGOPlus²¹ were introduced after the last
256 CAFA competition, and they were shown to perform very well against top CAFA algorithms.
257 (For details of the implementation, parameter settings and a description of these algorithms
258 see Materials and Methods and Supplementary Notes 14, 16 and 17).

259

260 Figures 2-5 show the AUC-ROC, AUC-PR, F_{\max} and S_{\min} evaluated *per-gene* and the *per-term*
261 for S2F and each competitor algorithm. (An interactive version of these results is also
262 available in the result explorer on our website: <https://www.paccanarolab.org/s2f>). We can
263 see that S2F outperformed the other methods according to the vast majority of the
264 performance measures for the ten bacteria – it is surpassed only in 4 out of the 80 bacteria-
265 measure combinations, most often on the AUC-ROC measure. In order to better appreciate
266 the increase in performance offered by S2F, we also explicitly report the percentage of
267 improvement of S2F vs each competitor for each of the 10 organisms (see Supplementary
268 Figures 53-59 in Supplementary Note 15).

269

270 Analysing these results, we see that, as expected, the accuracy of the S2F predictions does
271 depend on the accuracy of InterPro and HMMER, that provide the initial seeds for the
272 diffusion process of S2F. An interesting question is whether the improved performance of
273 S2F is merely due to the fact that it combines the labels of InterPro and HMMER, or whether
274 the diffusion of these labels through the transferred networks has a role in its performance.
275 For this reason, we also report in the figures the performance of the linear combination of
276 InterPro and HMMER labels that we used as seeds for the diffusion process in S2F (matrix
277 Y). We can see that, with the only exception of the AUC-ROC for *Brucella Abortus*, S2F

278 shows an improvement when compared with the simple linear combination of the InterPro
279 and HMMER outputs. This means that S2F is able to effectively combine the information of
280 these labels together with the evolutionary information contained in the interolog graphs.

281

282 As we mentioned earlier, by integrating InterPro and HMMER we aimed at obtaining seeds
283 that combined the high accuracy and specificity offered by InterPro with the high coverage
284 provided by HMMER. To check whether our linear combination, controlled by the
285 parameter α achieved this, we analysed how the different setting of α affected the S2F
286 results (details of the experiments are described in Supplementary Note 13). Supplementary
287 Figures 48-51 show that, in general, a combination of InterPro and HMMER seeds
288 ($0 < \alpha < 1$) gives much better results in terms of S2F performance than when using only
289 seeds from either of them ($\alpha = 0$ or $\alpha = 1$). However, just looking at S2F performance, it is
290 unclear how to set the value of α , as there is disagreement among different performance
291 measures and organisms. At the same time, an important objective in real-world scenarios
292 is to predict, for a given gene, a small set of terms that are highly accurate while being as
293 specific as possible. Therefore, we analysed the information content of the top genes
294 predicted by S2F for different values of α (see Supplementary Figure 52). Our results show
295 that, in this scenario, high values of α (e.g., $\alpha = 0.9$) should be preferred.

296

297 We also evaluated the predictions obtained by diffusing the outputs of InterPro and
298 HMMER, separately, on the interolog network W . Supplementary Figures 11-14
299 (Supplementary Note 8, also available in the interactive data explorer on our website
300 <https://www.paccanarolab.org/s2f>) show how our diffusion process is able to improve the
301 labels obtained by InterPro (or HMMER). This means that our diffusion on combined

302 interolog networks is an effective way to improve protein function prediction over simpler
303 homology methods.

304

305 Our diffusion method was motivated by our desire to model the presence of overlapping
306 communities in functional networks. It is unclear how to quantify exactly the number of
307 proteins being shared across communities, as this is obscured by the relationships among
308 functional labels as well as the noise and incompleteness of available annotations. However,
309 the semantic similarity of proteins with known function can provide some insight, as we can
310 quantify the correlation between the graph onto which we diffuse, W^{S2F} , and a graph of
311 semantic similarities among functionally annotated proteins, G^{SS} . Supplementary Figure 17
312 shows the values of these correlations for each of the ten bacteria and compares them with
313 correlations between G^{SS} and W^{GM} , the graph used by GeneMANIA¹⁷, a diffusion-based
314 method for protein function prediction in model organisms that does not explicitly model
315 overlapping communities (for details of these experiments see Supplementary Note 6). We
316 can see that W^{S2F} shows higher correlation with the semantic similarity graph G^{SS} in the
317 great majority of the cases, for different organisms and across different GO ontologies.

318

319 Finally, to further demonstrate how S2F can facilitate biological research by generating
320 feasible hypothesis, we performed a prospective evaluation. We deployed S2F to make
321 predictions using only data available up to December 2014 and we assessed its accuracy on
322 proteins that were experimentally annotated between 2015 and 2021. The experiments are
323 detailed in Supplementary Note 11. Supplementary Figures 44-47 show that while the
324 performance of InterPro is relatively stable, for some bacteria the overall performance of
325 HMMER (and, as a consequence, of the InterPro + HMMER combination) seems to worsen
326 greatly. As expected, the performance of S2F decreases in these cases, but overall the

327 diffusion process is able to alleviate the effect and compensate for the lower quality of the

328 seeds.

329

330 **Discussion**

331 The difficulty of protein function prediction, one of the most important problems in
332 computational biology, varies greatly, depending on how much experimental information is
333 available for the organism under investigation. Predictions for well-studied organisms can
334 rely on multiple types of functional experimental evidence (e.g., gene expression,
335 proteomics data) that can be represented in the form of graphs. For these organisms,
336 network propagation approaches that amplify existing knowledge about the function of
337 some of the proteins have been shown to be very effective^{5,17,24,25}.

338

339 This paper introduces S2F, a method that applies a network propagation algorithm to
340 organisms for which only sequence information is available. The main idea is to create
341 networks of interologs by systematically transferring functional data that is available for
342 model organisms, and to use these networks to combine and amplify a few preliminary GO
343 labels (seeds) obtained through homology or identifiable protein features.

344

345 Our work shows that employing a diffusion process over networks of interologs is an
346 effective way to improve predictions over simple homology. The improvement comes from
347 combining information: S2F effectively integrates homology information and identifiable
348 protein features (preliminary GO labels from HMMER and InterPro) together with
349 evolutionary information contained in the interolog graphs, through a diffusion process. S2F
350 includes a novel network propagation algorithm that can account for the presence of
351 overlapping communities of nodes with related functions.

352

353 Ultimately, the accuracy of S2F when predicting function for a specific organism will depend
354 on several factors, including the specificity and diversity of the preliminary GO labels, and
355 the density of the interolog networks, which in turn depends on the evolutionary distance
356 from organisms with existing functional experimental evidence. When predicting a GO term
357 for a specific gene, these factors affect how many neighbours that gene has, how many of
358 these genes have preliminary GO labels, and how accurate these labels are. These factors
359 are highly interleaved, and it is difficult to quantify the effect of each one individually. For
360 example, it would seem reasonable to expect that S2F would generate better predictions for
361 more highly connected nodes. We tested this hypothesis by measuring the correlation
362 between node degrees and the performance measures for the bacteria in this study.
363 However, our results show that the correlation was either weak and negative, or not
364 statistically significant (see Supplementary Note 10, as well as Supplementary Figures 26-
365 36).

366

367 The different interolog networks that we combine are extremely sparse with virtually no
368 overlap among them (see Supplementary Figures 3 and 4). In this scenario, in terms of
369 prediction performance, different combination methods would give results that are as good
370 as the simple average of the networks (Supplementary Figures 5-8 compare our
371 combination strategy, the network combination used by STRING¹⁶, and the simple average).
372 However, our approach allows the linear combination of the different networks through
373 learned coefficients, which provides us with information about the relative importance and
374 role of each network in the prediction (see Supplementary Note 4). Our combination
375 method is similar to the one used in GeneMANIA, but it allows us to learn these linear
376 weights without relying on an initial set of known functional labels.

377

378 We note that the removal of functional information regarding each bacterium and its
379 phylogenetically close species, makes this problem much harder than the one tested in the
380 regular CAFA competition settings. For this reason, the performances for Argot 2.5²⁰,
381 DeepGOPlus²¹, GOLabeler²², and NetGO²³ seem generally lower than those reported earlier.
382 Also, methods that are able to integrate global and local information seem to perform
383 better than local methods in our setting. This can be seen by comparing the results obtained
384 by S2F and the Consistency Method²⁶ – another method that integrates global information –
385 with the results obtained by NetGO, where the use of network information is limited locally
386 to nodes that are just one link away from the query node. A performance comparison
387 between our label propagation method and the Consistency Method is available in
388 Supplementary Note 9.

389

390 In this paper we have focused and presented results for bacteria, but S2F can be applied to
391 any organism, independently on how well functionally characterised it is. An earlier version
392 of S2F which is optimised to use existing functional evidence for target organisms was
393 submitted to the CAFA2 challenge³, where it ranked as one of the top performing methods.

394

395 The code for S2F is freely available at <https://www.paccanarolab.org/s2f>. The S2F software
396 is fast, robust, and easy to setup and run. The software is fully documented, including a wiki
397 with instructions for common use cases, instructions on how to use S2F to predict function
398 for newly sequenced bacteria and details on how to replicate all our results, together with
399 the necessary input data (see Supplementary Data).

400 **Materials and Methods**

401 **S2F Seed Inference**

402

403 InterPro produces m binary matrices of predictions $R^{(k)}$, each of size $(n \times t)$ (here
404 $k \in \{1, \dots, m\}$ and $m \leq 14$ is the number of models for which InterPro gives at least one
405 prediction for the target organism). To combine these matrices while ensuring that the
406 combination is consistent with the hierarchical structure of GO, we first up-propagate these
407 associations according to the true path rule^{7,8} considering both the “is_a”, and “part_of”
408 relations. Each matrix $R^{(k)}$ is up-propagated separately, and therefore any convex
409 combination of the up-propagated matrices will be consistent. We combine them to obtain
410 a consistent InterPro seed set $R \in \mathbb{R}^{n \times t}$ where each entry of R , R_{ij} , is defined as:

411

$$R_{ij} = \frac{\sum_{k=1}^m R_{ij}^k}{m}$$

412

413 **S2F Network Transfer**

414

415 STRING¹⁶ is a database that compiles several 3,123,056,667 interactions between proteins in
416 5,090 organisms. Interactions are divided into 7 types: “neighborhood”, “fusion”, “co-
417 occurrence”, “experiments”, “co-expression”, “textmining”, and “database”. Each
418 interaction is annotated with a score that ranges from 0 to 1 representing the confidence
419 that STRING assigns for the two proteins to be functionally related.

420

421 In our transfer procedure, two proteins A and A' are considered to be orthologues if three
422 conditions are met:

- 423 • they are BLAST mutual best hits, with both e-values smaller than 1e-6;
- 424 • percent identity is greater than 80% – this is to avoid transference between multi-
425 domain proteins with different domain architecture;
- 426 • their “joint identity” (geometric mean of the two percent identities) is above 60% –
427 Yu et al.¹³ showed that this condition achieves almost perfect accuracy at identifying
428 interacting orthologues.

429

430 When the same interaction can be transferred from multiple organisms, only the one with
431 the highest “joint identity” is kept. The pseudocode of the algorithm for building a collection
432 of transferred networks for the target organism is provided in Supplementary Algorithm 1
433 (Supplementary Note 3). S2F only considers networks with at least 3 edges, that is, for every
434 interaction type in STRING, we consider the transferred network r only if $W^{(r)}$ contains at
435 least 3 values.

436

437 Finally, a homology network is added to the collection of interolog networks to increase the
438 combined network connectivity and facilitate the diffusion process. The homology network
439 $W^{(h)}$ is defined as the negative log of the BLAST e-value for every pair of proteins.

440

441 **S2F Network combination**

442

443 Given p networks $W^{(r)}$ with $r \in \{1, \dots, p\}$, we combine them into a single network W using
 444 a weighted linear combination. The vector of weights $\hat{c} \in \mathbb{R}^p$, and bias \hat{b} are learnt by
 445 minimising:

$$(\hat{c}, \hat{b}) = \underset{c, b}{\operatorname{argmin}} \sum_{i, j} \left(b + \sum_{r=1}^p c_r W_{ij}^{(r)} - T_{ij} \right)^2$$

447

448 This linear regression can be solved efficiently, and we can interpret each learnt
 449 coefficient c_r as representing how much each network r contributes to the combination. An
 450 analysis on these coefficients is provided in the Supplementary Note 4.

451

452 **S2F label propagation**

453

454 The weighted Jaccard coefficient matrix J is defined elementwise as:

455

$$J_{ij} = \frac{\sum_k W_{ik} W_{jk}}{\sum_k W_{ik} + \sum_k W_{jk} - \sum_k W_{ik} W_{jk}}$$

456

457

458 Thus, the element J_{ij} relates to how much elements i and j belong to the same community
 459 in network W . For a given term k , we learn the k -th column of matrix F , that we denote by
 460 F_k , by minimising the cost function $Q(F_k)$:

461

$$Q(F_k) = \sum_{i=1}^n (F_{ik} - Y_{ik})^2 + \frac{\lambda}{2} \sum_{i=1}^n \frac{1}{d_i} \sum_{j=1}^n J_{ij} W_{ij} (F_{ik} - F_{jk})^2$$

462 Similarly to the cost function used by the Consistency Method (CM)²⁶, ours is the sum of two
 463 terms. The role of the first term is to conserve the initial labels Y_{ik} – this term is minimised
 464 when the node labels F_{ik} are the same as the initial labels. The second term accounts for
 465 the consistency of the labels of adjacent nodes (reflecting the guilt-by-association principle)
 466 – this term is minimised when adjacent nodes have similar labels (i.e. the difference
 467 between F_{ik} and F_{jk} becomes small). Note that the importance of the difference between
 468 F_{ik} and F_{jk} is proportional to $J_{ij}W_{ij}$ which models the community effect — the more i and j
 469 are connected through their neighbours, the greater their contribution to the cost function.
 470 Furthermore, notice that:

471

$$\frac{1}{d_i} = \frac{1}{\sum_j J_{ij}W_{ij}}$$

472

473 is a normalisation factor that gives to each protein in the network similar ability to influence
 474 its neighbours, independently of its degree.

475 The closed-form solution that minimises $Q(F_k)$ is:

476

$$F_k^* = (I + \lambda L)^{-1}Y_k$$

477

478 where Y_k is the initial labelling, $L = D^{S2F} - W^{S2F}$ is the Laplacian of W^{S2F} , whose entry
 479 W_{ij}^{S2F} is defined as:

480

481

$$W_{ij}^{S2F} = \frac{1}{2} \left(\frac{1}{d_i} + \frac{1}{d_j} \right) J_{ij}W_{ij},$$

482

483 and D^{S2F} is a diagonal matrix where the i -th diagonal element is $D_{ii}^{S2F} = \sum_j W_{ij}^{S2F}$.

484 **Bacteria selection criteria and Datasets**

485 The criteria we used for selecting bacteria were:

- 486 • The bacteria must have at least 10 functional annotations with an experimental or
487 curated GO evidence code (EXP, IDA, IPI, IMP, IGI, IEP, TAS, or IC) in the GOA
488 database¹⁹.
- 489 • The bacteria must have at least 8 terms annotated with at least 3 genes after up-
490 propagation, for each GO subdomain – biological process (BP), molecular function
491 (MF), and cellular component (CC).

492

493 In our experiments, we used STRING version 11.0. All sequences in FASTA format were
494 downloaded from UniProtKB/Swiss-Prot using the taxonomy identifiers listed in Table 1. The
495 GO annotations were downloaded from the GOA database¹⁹. All datasets were downloaded
496 in April 2020. We used HMMER version 3.1b2, InterProScan version 5.42-78.0, and blastp
497 from BLAST 2.6.0+.

498

499 **Competitor algorithms**

500

501 In all our experiments, in order to simulate a real case scenario for the problem of
502 predicting function in newly sequenced organisms, for each bacterium, we removed any
503 functional information regarding that bacterium as well as any functional information about
504 species that are phylogenetically close (the list of all organisms excluded is provided in the
505 Supplementary Data).

506

507 GOLabeler²² and its successor, NetGO²³, are only offered as web services which use all the
508 data available from their sources (namely GOA, STRING, UniProtKB, InterPro) for their
509 prediction. Therefore, the results for NetGO and GOLabeler presented here were obtained
510 running our own implementation of these systems that had been trained using datasets
511 from which all the aforementioned functional information had been removed. All the
512 parameters of the component models as well as the learning to rank ensemble were set
513 using the default values suggested by the authors^{22,23}. A detailed description on how to
514 prepare the input data and how to use our implementation of these methods is available in
515 Supplementary Note 17.

516

517 Argot 2.5²⁰ was run on its web server (<http://www.medcomp.medicina.unipd.it/Argot2-5/>).
518 For each bacterium, we first used BLAST and HMMER to obtain alignments between its
519 proteins and a version of UniProtKB from which the sequences of excluded organisms (for
520 that bacterium) were omitted. These alignments were then submitted to the Argot 2.5 web
521 server.

522

523 DeepGOPlus²¹ was run using the code from the latest stable version available (1.0.1). In
524 order to remove the information from phylogenetically close organisms, we added some
525 pre-processing steps to the input files and small corrections were made to the prediction
526 script. A detailed guide on how to setup and run the pre-processing and prediction is
527 described in Supplementary Note 16.

528

529 InterPro was run using InterProScan version 5.42-78.0, the output file was then processed to
530 extract the predictions that included GO terms.

531

532 HMMER version 3.1b2 was run against a GO annotation file that was pre-processed to keep
533 only the experimental or curated evidence codes (EXP, IDA, IPI, IMP, IGI, IEP, TAS, or IC). The
534 output file was post-processed to remove any alignment that came from an organism that
535 had been excluded in the prediction.

536 **Acknowledgements**

537 The first idea for this project was conceived in discussions with Tara Gianoulis. We
538 remember Tara dearly for her intelligence, kindness, enthusiasm and passion for research.
539 The authors also thank Prajwal Bhat, Tamás Nepusz, Juan Caceres, Marco Frasca, Giorgio
540 Valentini, Alessandra Devoto, Laszlo Bögre, Rajkumar Sasidharan, and Mark Gerstein for
541 many important and stimulating discussions.

542 A.P. was supported by Biotechnology and Biological Sciences Research Council
543 (<https://bbsrc.ukri.org/>) grants BB/K004131/1, BB/F00964X/1 and BB/M025047/1; Medical
544 Research Council (<https://mrc.ukri.uk>) grant MR/T001070/1; Consejo Nacional de Ciencia y
545 Tecnología Paraguay (<https://www.conacyt.gov.py/>) grants 14-INV-088 and PINV15-315;
546 National Science Foundation Advances in Bio Informatics (<https://www.nsf.gov/>) grant
547 1660648; Fundação de Amparo à Pesquisa do Estado do Rio de Janeiro grant E-
548 26/201.079/2021 (260380); and Fundação Getulio Vargas.

549 **Author Contributions**

550 A.P. conceived the study. A.P. and H.Y. devised the algorithms, developed the prototype and
551 performed preliminary evaluations. M.T. and A.E.R. implemented and extended the
552 algorithms and evaluation metrics, performed large scale experiments and analysed the
553 results. A.P., M.T., A.E.R. wrote the manuscript and evaluated the biological relevance of the
554 results. All the authors discussed the results and implications. A.P. supervised the project.

555 **Competing Interests**

556

557 The authors declare no competing interests.

558

559 **Data Availability**

560

561 The input sequence files²⁷ in FASTA format for all the organisms used in this paper, are

562 available at: <https://doi.org/10.5281/zenodo.5514323>. The same URL also contains the

563 detailed list of all organisms excluded when testing each specific bacterium.

564

565 **Code Availability**

566

567 The code for S2F is freely available and maintained at <https://www.paccanarolab.org/s2f>.

568 The exact version²⁸ used for this publication is available at:

569 <https://doi.org/10.5281/zenodo.5513071>

570 **Tables**

571

572 **Table 1** List of bacteria that satisfies our selection criteria with number of genes and

573 annotations. The number of terms with more than 3 annotations in each of the GO domains

574 (BP = biological process, MF = molecular function, CC = cellular component) is calculated

575 after up-propagation, and therefore may be larger than the number of experimentally

576 annotated genes.

NCBI ID	Name	Genes	Experimentally annotated genes	BP terms with > 3 annotations	MF terms with > 3 annotations	CC terms with > 3 annotations
272624	<i>Legionella pneumophila</i> subsp. <i>Pneumophila Philadelphia 1</i>	2,076	18	30	8	8
223283	<i>Pseudomonas syringae</i> pv. <i>tomato</i>	5,055	25	48	32	15
359391	<i>Brucella abortus</i>	2,229	26	17	8	14
99287	<i>Salmonella typhimurium</i>	3,764	116	183	46	24
198628	<i>Dickeya dadantii</i>	3,411	102	214	21	13
1111708	<i>Synechocystis</i> sp.	2,442	137	101	21	30
224308	<i>Bacillus subtilis</i>	3,410	375	301	120	24
208964	<i>Pseudomonas aeruginosa</i>	4,487	947	695	222	42
83332	<i>Mycobacterium tuberculosis</i>	3,284	1,027	797	280	45
83333	<i>Escherichia coli</i>	3,906	3,350	1,546	706	134

577

579 **Figure Legends/Captions**

580

581 **Figure 1 Overview of the S2F approach** The set of n protein sequences of the target
582 organism (shown in red) constitutes the input to the system; t is the total number of GO
583 terms to be predicted. External datasets (STRING ¹⁶, GOA ¹⁹, and UniProtKB ²⁹) are shown in
584 orange. **Seed Inference.** Running HMMER on the input sequences against experimentally
585 annotated sequences from UniProtKB/Swiss-Prot we obtain an $(n \times t)$ matrix H of
586 predictions (the HMMER seed set). Running InterPro we obtain m matrices of predictions
587 $R^{(m)}$, one per InterPro model, each of size $(n \times t)$. These matrices are then combined into a
588 single $(n \times t)$ matrix R (the InterPro seed set). The combined seed set Y , that will be used
589 for the label propagation, is a linear combination of H and R . **Network Transfer.** A collection
590 of networks is built by our interaction transfer procedure using known functional
591 relationships between proteins in every organism from the STRING database. **Network**
592 **Combination.** Transferred networks are linearly combined into a single network W . The
593 weights of the linear combination are learnt using an auxiliary target network built from R .
594 **Prediction.** The network W and the seed set Y are fed into our label propagation algorithm
595 that outputs the protein function prediction F , an $(n \times t)$ matrix where each row
596 corresponds to a protein, each column corresponds to a GO term and each entry $F_{i,j}$ is
597 related to the probability for protein i to have function j . For a given protein i , its labels $F_{i,\cdot}$
598 are guaranteed to be consistent, i.e., they satisfy the GO “true path rule”.

599

600 **Figure 2** S_{\min} metric for every organism per-gene (left) and per-term (right), lower values are
601 better. Comparison of HMMER, InterPro, HMMER + InterPro, S2F, Argot 2.5, DeepGOPlus,
602 GOLabeler, and NetGO. Values indicate the mean of the metric over genes or terms, and
603 error bars indicate a confidence interval of 95%, estimated using 10,000 bootstrap iterations
604 on the gene set or term set, respectively.

605

606 **Figure 3** F_{\max} for every organism per-gene (left) and per-term (right), higher values are
607 better. Comparison of HMMER, InterPro, HMMER + InterPro, S2F, Argot 2.5, DeepGOPlus,
608 GOLabeler, and NetGO. Values indicate the mean of the metric over genes or terms, and
609 error bars indicate a confidence interval of 95%, estimated using 10,000 bootstrap iterations
610 on the gene set or term set, respectively.

611

612 **Figure 4** AUC-ROC for every organism per-gene (left) and per-term (right), higher values are
613 better. Comparison of HMMER, InterPro, HMMER + InterPro, S2F, Argot 2.5, DeepGOPlus,
614 GOLabeler, and NetGO. Values indicate the mean of the metric over genes or terms, and
615 error bars indicate a confidence interval of 95%, estimated using 10,000 bootstrap iterations
616 on the gene set or term set, respectively.

617

618 **Figure 5** AUC-PR for every organism per-gene (left) and per-term (right), higher values are
619 better. Comparison of HMMER, InterPro, HMMER + InterPro, S2F, Argot 2.5, DeepGOPlus,
620 GOLabeler, and NetGO. Values indicate the mean of the metric over genes or terms, and
621 error bars indicate a confidence interval of 95%, estimated using 10,000 bootstrap iterations
622 on the gene set or term set, respectively.

623

624 **References**

625

- 626 1. Cruz, L. M., Trefflich, S., Weiss, V. A. & Castro, M. A. A. Protein Function Prediction.
627 *Methods Mol. Biol. Clifton NJ* **1654**, 55–75 (2017).
- 628 2. Shehu, A., Barbará, D. & Molloy, K. A Survey of Computational Methods for Protein
629 Function Prediction. in *Big Data Analytics in Genomics* (ed. Wong, K.-C.) 225–298
630 (Springer International Publishing, 2016). doi:10.1007/978-3-319-41279-5_7.
- 631 3. Jiang, Y. *et al.* An expanded evaluation of protein function prediction methods shows an
632 improvement in accuracy. *Genome Biol.* **17**, 184 (2016).
- 633 4. Ashburner, M. *et al.* Gene Ontology: tool for the unification of biology. *Nat. Genet.* **25**,
634 25–29 (2000).
- 635 5. Cowen, L., Ideker, T., Raphael, B. J. & Sharan, R. Network propagation: a universal
636 amplifier of genetic associations. *Nat. Rev. Genet.* **18**, 551–562 (2017).
- 637 6. Zhou, N. *et al.* The CAFA challenge reports improved protein function prediction and
638 new functional annotations for hundreds of genes through experimental screens.
639 *Genome Biol.* **20**, 244 (2019).
- 640 7. Valentini, G. True Path Rule Hierarchical Ensembles for Genome-Wide Gene Function
641 Prediction. *IEEE/ACM Trans. Comput. Biol. Bioinform.* **8**, 832–847 (2011).
- 642 8. Friedberg, I. & Radivojac, P. Community-Wide Evaluation of Computational Function
643 Prediction. in *The Gene Ontology Handbook* (eds. Dessimoz, C. & Škunca, N.) 133–146
644 (Springer New York, 2017). doi:10.1007/978-1-4939-3743-1_10.
- 645 9. Obozinski, G., Lanckriet, G., Grant, C., Jordan, M. I. & Noble, W. S. Consistent
646 probabilistic outputs for protein function prediction. *Genome Biol.* **9**, S6 (2008).
- 647 10. Mitchell, A. L. *et al.* InterPro in 2019: improving coverage, classification and access to
648 protein sequence annotations. *Nucleic Acids Res.* **47**, D351–D360 (2019).

- 649 11. HMMER. <http://hmmer.org/>.
- 650 12. Walhout, A. J. *et al.* Protein interaction mapping in *C. elegans* using proteins involved in
651 vulval development. *Science* **287**, 116–122 (2000).
- 652 13. Yu, H. *et al.* Annotation Transfer Between Genomes: Protein–Protein Interologs and
653 Protein–DNA Regulogs. *Genome Res.* **14**, 1107–1118 (2004).
- 654 14. Ben-Hur, A. & Noble, W. S. Kernel methods for predicting protein–protein interactions.
655 *Bioinformatics* **21**, i38–i46 (2005).
- 656 15. Sharan, R. *et al.* Conserved patterns of protein interaction in multiple species. *Proc. Natl.*
657 *Acad. Sci.* **102**, 1974–1979 (2005).
- 658 16. Szklarczyk, D. *et al.* STRING v11: protein–protein association networks with increased
659 coverage, supporting functional discovery in genome-wide experimental datasets.
660 *Nucleic Acids Res.* **47**, D607–D613 (2019).
- 661 17. Mostafavi, S., Ray, D., Warde-Farley, D., Grouios, C. & Morris, Q. GeneMANIA: a real-
662 time multiple association network integration algorithm for predicting gene function.
663 *Genome Biol.* **9**, S4 (2008).
- 664 18. Zhou, N. *et al.* The CAFA challenge reports improved protein function prediction and
665 new functional annotations for hundreds of genes through experimental screens.
666 *bioRxiv* 653105 (2019) doi:10.1101/653105.
- 667 19. Huntley, R. P. *et al.* The GOA database: gene Ontology annotation updates for 2015.
668 *Nucleic Acids Res.* **43**, D1057–1063 (2015).
- 669 20. Lavezzo, E., Falda, M., Fontana, P., Bianco, L. & Toppo, S. Enhancing protein function
670 prediction with taxonomic constraints – The Argot2.5 web server. *Methods* **93**, 15–23
671 (2016).
- 672 21. Kulmanov, M. & Hoehndorf, R. DeepGOPlus: improved protein function prediction from
673 sequence. *Bioinformatics* **36**, 422–429 (2020).

- 674 22. You, R. *et al.* GOLabeler: improving sequence-based large-scale protein function
675 prediction by learning to rank. *Bioinformatics* **34**, 2465–2473 (2018).
- 676 23. You, R. *et al.* NetGO: improving large-scale protein function prediction with massive
677 network information. *Nucleic Acids Res.* **47**, W379–W387 (2019).
- 678 24. Makrodimitis, S., van Ham, R. C. H. J. & Reinders, M. J. T. Automatic Gene Function
679 Prediction in the 2020's. *Genes* **11**, 1264 (2020).
- 680 25. Cao, M. *et al.* Going the Distance for Protein Function Prediction: A New Distance Metric
681 for Protein Interaction Networks. *PLOS ONE* **8**, e76339 (2013).
- 682 26. Zhou, D., Bousquet, O., Lal, T. N., Weston, J. & Schölkopf, B. Learning with Local and
683 Global Consistency. in 8 (2004).
- 684 27. Torres, M., Haixuan Yang, Romero, A. E. & Paccanaro, A. Input Data for 'Protein Function
685 Prediction for newly sequenced organisms'. (2021) doi:10.5281/ZENODO.5514323.
- 686 28. Torres, M., Haixuan Yang, Romero, A. E. & Paccanaro, A. *Source code for 'Protein
687 Function Prediction for newly sequenced organisms'*. (Zenodo, 2021).
688 doi:10.5281/ZENODO.5513071.
- 689 29. UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res.* **47**, D506–D515
690 (2019).
- 691









