

Biases in the Experimental Annotations of Protein Function and their Effect on Our Understanding of Protein Function Space

1 Alexandra Schnoes Dept/Program/Center, University of California, San Francisco, San Francisco, CA, USA

2 David C. Ream Department of Microbiology, Miami University, Oxford, OH, USA

3 Alexander W. Thorman, Department of Microbiology, Miami University, Oxford, OH, USA

4 Patricia C. Babbitt Dept/Program/Center, University of California, San Francisco, San Francisco, CA, USA

5 Iddo Friedberg, Department of Microbiology and Computer Science & Software engineering, Miami University, Oxford, OH, USA

* E-mail: corresponding i.friedberg@muohio.edu

1 Abstract

2 **Background:** Computational protein function prediction programs rely upon well-annotated
3 databases for testing and training their algorithms. These databases, in turn, rely upon
4 the work of curators to capture experimental findings from scientific literature and ap-
5 ply them to protein sequence data. However, with the increasing use of high-throughput
6 experimental assays, a small number of experimental articles dominate the functional
7 protein annotations collected in databases. Here we investigate just how prevalent is the
8 “few articles – many proteins” phenomenon. We hypothesize that the dominance of high-
9 throughput experiments in proteins annotation biases our view of the corpus of functions

10 enabled by proteins.

11 **Results:** We examine the annotation of UniProtKB by the Gene Ontology Annotation
12 project (GOA), and show that the distribution of proteins per article is exponential, with
13 0.06% of articles dominating 20% of the annotations. Since each of the dominant articles
14 describes the use of an assay that can find only one function or a small group of functions,
15 this leads to substantial biases, in several aspects, in what we know about the function of
16 many proteins.

17 **Conclusions:** Given the experimental techniques available, protein function annota-
18 tion bias due to high-throughput experiments is unavoidable. Knowing that these biases
19 exist and understanding their characteristics and extent is important for database cura-
20 tors, developers of function annotation programs, and anyone who uses protein function
21 annotation data to plan experiments.

22 Author Summary

23 Introduction

24 Functional annotation of proteins is an open problem and a primary challenge in molec-
25 ular biology today [1–4]. The ongoing improvements in sequencing technology had the
26 emphasis shifting from realizing the \$1000 genome to the 1-hour genome [5]. The ability
27 to rapidly and cheaply sequence genomes is creating a flood of sequence data, but to
28 make these data useful, extensive analysis is needed. A large proportion of this work
29 involves assigning biological function to newly determined gene sequences, a process that
30 is both complex and costly [6]. To aid current annotation procedures and improve com-

putational function prediction algorithms, sources of high-quality, experimentally derived functional data are necessary. Currently, one of the few repositories of such data is the UniProt-GOA database [7], which contains both computationally derived and literature derived functional information. The literature derived information is extracted by human curators who capture functional data from publications, assign the data to its appropriate place in the Gene Ontology hierarchy [8] and label them with appropriate functional evidence codes. The UniProt-GOA database is one of only a small number of databases that explicitly connects functional data, publication references and evidence codes to experimentally studied proteins. In addition, annotations captured in UniProt-GOA directly impact the annotations in the UniProt/Swiss-Prot database, widely considered to be a gold standard set of functional annotation [2]. It is therefore important to understand any trends and biases that are encapsulated by the UniProt-GOA database, as those impact well-used sister databases and therefore a large number of users worldwide. Furthermore, any biases would impact function prediction algorithms development and training.

One concern surrounding the capture of functional data from articles is the propensity for high-throughput experimental work to become a large fraction of the data in UniProt-GOA, thus having a small number of experiments dominate the protein function landscape. In this work we analyzed the relative contribution of peer-reviewed articles describing all the experimentally-derived annotations in UniProt-GOA. We found some striking biases, stemming from the fact that a small fraction of articles describing high-throughput experiments disproportionately contribute to the pool of experimental annotations of model organisms. Consequently, we show that: 1) annotations coming from high-throughput experiments are overall less informative than those provided by low-throughput experiments; 2) annotations from high throughput experiments bias the annotations towards a limited number of functions, and, 3) many high-throughput ex-

periments overlap in the proteins they annotate, and in the annotations assigned. Taken together, our findings offer a picture of how the protein function annotation landscape is generated from scientific literature. Furthermore, due to the biases inherent in the current system of sequence annotations, this study serves as a caution to the producers and consumers of biological data from high-throughput experiments.

Methods and Results

Articles and Proteins

With the advent of high-throughput experiments it has become possible to conduct large-scale studies of protein functions. Consequently, some studies reveal certain functional aspects of a large amount of proteins as a result of the particular type of assay or assays used. To understand the impact of large-scale studies on the corpus of experimentally annotated proteins, we looked at the UniprotKB Gene Ontology (GO) Annotation database, or UniProt-GOA. Proteins in UniProt-GOA have been annotated with one or more GO terms using a procedure described in [7]. Briefly, this procedure consists of six steps which include sequence curation, sequence motif analyses, literature-based curation, reciprocal BLAST [9] searches, attribution of all resources leading to the included findings, and quality assurance. If the annotation source is a research article, the attribution includes its PubMed ID. For each GO term associated with a protein, there is also an *evidence code* (EC) which the curator assigns to explain how the association between the protein and the GO term was made. Experimental evidence codes include such terms as: Inferred by Direct Assay (IDA) which indicates that “a direct assay was carried out to determine the function, process, or component indicated by the GO term” or Inferred from Physical Interaction (IPI) which “Covers physical interactions between the gene product of

interest and another molecule.” (All EC definitions were taken from the GO site, geneontology.org). Computational evidence codes include terms such as *Inferred from Sequence or Structural Similarity* (ISS) and *Inferred from Sequence Orthology* (ISO). Although the evidence in computational evidence codes is non-experimental, the proteins annotated with these evidence codes are still assigned by a curator, rendering a degree of human oversight. Finally, there are also computational, non-experimental evidence codes, the most prevalent being *Inferred from Electronic Annotation* (IEA) which is “used for annotations that depend directly on computation or automated transfer of annotations from a database”. IEA evidence means that the annotation is wholly automated, and was not made or checked by a person. Different degrees of reliability are associated with different evidence codes, with experimental codes generally considered to be of higher reliability than non-experimental codes. However, the increase in the number of high-throughput experiments used to determine protein functions may introduce biases into experimental protein annotations, due to the inherent capabilities and limitations of high-throughput assays. To test the hypothesis that such biases exist, and to study their extent if they do, we compiled the details of all experimentally-annotated proteins in UniProtKB. This included all proteins whose GO annotations have the GO experimental evidence codes EXP, IDA, IPI, IMP, IGI, IEP. We first examined the distribution of articles by the number of proteins they annotate. As can be seen in Figure 1, the distribution of the number of proteins annotated per article follows a power-law distribution. $f(x) = ax^k$. Using linear regression over the log values of the axes we obtained a fit with $p < 1.18 \times 10^{-8}$ and $R^2 = -0.72$. We therefore conclude that there is indeed a substantial bias in experimental annotations, in which there are few articles that annotate a large number of proteins.

To better understand the consequences of such a distribution, we divided the annotating articles into four cohorts, based on the number of proteins each article annotates.

104 *Single-throughput* articles are those articles that annotate only one protein; *low through-*
 105 *put* articles annotate 2-9 proteins; *moderate throughput* articles annotate 10-99 proteins
 106 and *high throughput* articles annotate over 99 proteins. The results are shown in Table 2.
 107 The most striking finding is that high throughput articles are responsible for 25% of the
 108 annotations in Uniprot-GOA, even though they comprise 0.08% of the articles. 96% of
 109 the articles are single-throughput and low throughput, however those annotate only 53%
 110 of the proteins in Uniprot-GOA. So while moderate throughput and high-throughput ex-
 111 periments account for almost half of the annotations in Uniprot-GOA, they comprise only
 112 4% of the experiments published.

113 What typifies high-throughput articles? Also, how may the log-odds distribution
 114 bias what we understand of the protein function universe? To answer these questions,
 115 we examined different aspects of the annotations in the four article cohorts. Also, we
 116 examined in higher detail the top 50 high-throughput annotating articles. (Overall, only
 117 108 articles in our study annotated more than 100 proteins). “Top-50 high throughput
 118 annotating articles” are those articles describing experimental annotations that are ranked
 119 by the number of of proteins annotated per article. An initial characterization of the top
 120 50 high-throughput articles is shown in Table. As can be seen, most of the articles are
 121 specific to a single species (typically a model organism) and to a single assaying pipeline
 122 that is used to assign function to the proteins in that organism. Typically only one
 123 ontology (MFO, BPO or CCO) was used for annotation. For some species this means
 124 that a single functional aspect (MFO, BPO or CCO) of a species will be dominated by a
 125 single study / publication.

126 Term frequency bias

127 To see how much a single species- and method-specific high-throughput assay affects the
 128 entire annotation of a species, we examined the relative contribution of the top-50 articles
 129 to the entire corpus of experimentally annotated protein in each species. Unsurprisingly,
 130 all the species found in the top-50 articles were either common model organisms or human.
 131 For each species, we looked at the five most frequent terms in the top 50 annotating
 132 articles. We then examined the contribution of this term by the top 50 articles to the
 133 general annotations of that species. The *contribution* is the number of annotations by
 134 any given GO term in the top 50 articles divided by the number of annotations by that
 135 GO term in all of UniProtKB. For example, as seen in Figure 3 in *D. melanogaster* 88%
 136 of the use of the term “precatalytic spliceosome” in all articles experimentally annotating
 137 this species is contributed by the top-50 articles.

138 For most organisms in the top-50 articles, the annotations were within the cellular
 139 component ontology. Notable exceptions are *D. melanogaster* and *C. elegans* where the
 140 dominant terms were from the Biological Process ontology, and in mouse, where “protein
 141 binding” and “identical protein binding” are from the Molecular Function Ontology. *D.*
 142 *melanogaster*’s annotation for the top terms is dominated (over 50% contribution) by the
 143 top-50 articles.

144 The term frequency bias described here can be viewed more broadly within the ontol-
 145 ogy bias. The proteins annotated by the cohorts of single-protein articles, low-throughput
 146 articles, and moderate throughput articles have similar ratios of the fraction of proteins
 147 annotated. Twenty-two to twenty-six percent of assigned terms are in the Molecular
 148 Function Ontology, and 51-57% are in the Biological Process Ontology and the remaining
 149 17-25% are in the Cellular Component ontology. These ratios change dramatically with
 150 high-throughput articles (over 99 terms per article). In the high-throughput articles, only

151 5% of assigned terms are in the Molecular Function Ontology, 38% in the Biological Pro-
 152 cess Ontology and 57% in the Cellular Compartment Ontology, ostensibly due to a lack of
 153 high-throughput assays that can be used for generating annotations using the Molecular
 154 Function Ontology.

155 Reannotation

156 Another type of annotation bias is due to re-annotation. How many of the top-50 articles
 157 actually re-annotate the same set of proteins? And how much of an agreement is there
 158 between different experiments? To investigate the extent of repetitive annotations in
 159 different articles, we clustered all the proteins annotated by the top-50 articles using CD-
 160 HIT [10] at 100% sequence identity. We then examined the number of clusters containing
 161 100% identical sequences per model species. The product of the number of proteins
 162 divided by the number of clusters is the redundancy percentage. For example, if each of
 163 the top-50 articles annotating the proteins in a given species annotated the same protein
 164 set, the redundancy percentage would be 100%. The results of the reannotation bias
 165 analysis are shown in Figure 2 and in Table 3. As can be seen, the highest redundancy
 166 (65%) is in the 12 articles annotating *C. elegans*.

167 We have determined therefore, that there is a varying degree of repetition between
 168 experiments in the proteins they annotate, with some overlaps being quite high. In those
 169 cases, many of the same proteins in the same organism are being annotated. However,
 170 there is still a need to determine whether this annotation is consistent or not. To do this,
 171 we looked for the proteins that are annotated by more than one article, within the same
 172 ontology.

173 Given a protein P , let G be the GO-terms g_1, g_2, \dots, g_m that annotate that protein
 174 in all top-50 articles for a single ontology $O \in \{BPO, MFO, CCO\}$. The count of each

PubMedID	UniProt ID	Ontology	GO-term	description
14562095	P36023	CCO	GO:0005634	nucleus
14562095	P36023	CCO	GO:0005737	cytoplasm
16823961	P36023	CCO	GO:0005739	mitochondrion
14576278	P36023	CCO	GO:0005739	mitochondrion

of these go terms per protein per ontology is n_1, n_2, \dots, n_m with n_i being the number of times GO term g_i annotates protein P .

The number of total annotations for a protein in an ontology is $\sum_1^m n_i$. The *maximum annotation consistency* for protein P in ontology O $0 \leq k_{P,O} \leq 1$ is calculated as:

$$k_{P,O} = \frac{\max(n_1, n_2, \dots, n_m)}{\sum_1^m n_i}; \text{formax}(n_1, n_2, \dots, n_m) \geq 2$$

For example, the protein “Oleate activated transcription factor 3” (UniProtID: P36023) in *S. cerevisiae* is annotated four times by three articles using the Cellular Component ontology:

The annotation consistency for P36023 is therefore the maximum count of identical GO terms (*mitochondrion*, 2), divided by the total number of annotations, 4: 0.5.

Table 4 shows the results of this analysis. In *A. thaliana*, 1941 proteins are annotated by 15 articles and 18 terms in the Cellular Component ontology. The mean maximum-consistency is 0.251. The highest mean consistency is for the annotation of 807 mouse proteins annotated in Cellular Component ontology with an annotation consistency 0.832. However, that is not surprising given that there are only three annotating articles, and two annotating terms. We omitted the ontology and organism combinations that were annotated by less than three articles or two GO terms, or both.

191 Quantifying annotation information

192 A common assumption holds that while high-throughput experiments do annotate more
 193 protein functions than low-throughput experiments, the former also tend to be more
 194 shallow in the predictions they provide. The information provided, for example, by a
 195 large-scale protein binding assay will only tell us if two proteins are binding, but will
 196 not reveal whether that binding is specific, will not provide an exact K_{bind} , will not say
 197 under what conditions binding takes place, or whether there is any enzymatic reaction
 198 or signal-transduction involved. Having on hand data from experiments with different
 199 “thoroughputness” levels, we set out to investigate whether there is a difference in the in-
 200 formation provided by high-throughput experiments vs. low-throughput ones. To answer
 201 this question, we first had to quantify the information given by GO terms. One way to do
 202 so, is to use the depth of the term in the ontology: the term “catalytic activity” (one edge
 203 distance from the ontology root node) would be less informative than “hydrolase activity”
 204 (two edges) and the latter will be less informative than “haloalkane dehalogenase activity”
 205 (five edges). We therefore counted edges from the ontology root term to the GO-term
 206 to determine term information. The larger the number of edges, the more specific –and
 207 therefore informative– is the annotation. In cases where several paths lead from the root
 208 to the examined GO-term, we used the minimal path. We did so for all the annotating
 209 articles split into groups by the number of proteins each article annotates.

210 Edge counting provides a measure of term-specificity. This measure is, however, im-
 211 perfect. The reason is that different areas of the GO DAG have different connectivities,
 212 and terms may have different depths unrelated to the intuitive specificity of a term.
 213 For example “D-glucose transmembrane transporter activity”, (GO:0055056) is 10 terms
 214 deep, while “L-tryptophan transmembrane transporter activity”, (GO:0015196) is four-
 215 teen terms deep. It is hard to discern whether these differences are meaningful. For this

reason, information content, the logarithm of the inverse of the GO term frequency in the corpus is generally accepted as a measure of GO-term information content [11, 12]. To account for the possible bias created by the GO-DAG structure, we also used the log-frequency of the terms in the experimentally annotated proteins in Uniprot-GOA. However, it should be noted that the log-frequency measure is also imperfect because, as we see throughout this study, a GO-term’s frequency may be heavily influenced by the top annotating articles, injecting a circularity problem into the use of this metric. Since no single metric for measuring the information conveyed by a GO term is wholly satisfactory, we used both edge-counting and information-content in this study.

The results of both analyses are shown in Figure 4. In general, the results from the depth-based analysis and the log-frequency based analysis are in agreement, when compared across groupings based on the number of proteins annotated by the articles. For the Molecular Function ontology, the distribution of edge counts and log-frequency scores decreases as the number of annotated proteins per-article increases. For the Biological Process ontology, the decrease is significant. However the contributor to the decrease are the high-throughput articles while there is little change in the first three article cohorts. Finally, there is no significant trend of GO-depth decrease in the Cellular Component Ontology. However, using the information content metric, there is also a significant decrease in information content in the high-throughput article cohort.

Evidence and Assertion

There are two complementary ways by which we come to know about a protein’s function. The twenty GO evidence codes, discussed above, encapsulate the type results by which the function was inferred, but they do not capture all the necessary information. For example, “Inferred by Direct Assay (IDA)” informs that experimental evidence was used, but does

not say which type of experiment was performed. This information is often needed, since knowing which experiments were performed can help the researcher establish the reliability and scope of the produced data. For example, RNA used in an RNAi experiment does not traverse the blood-brain-barrier, meaning that no data from the central nervous system can be drawn from an RNAi experiment. The Evidence Code Ontology, or ECO, seeks to improve upon the GO-attached evidence codes. ECO provides more elaborate terms than “Inferred by Direct Assay”: ECO also conveys which assay was used, e.g. “microscopy”, “RNA interference”. In addition to evidence terms, the ECO ontology provides *assertion terms* in which the nature of the assay is given. For example, an enzyme-linked immunosorbent assay (ELISA) provides quantitative protein data *in vitro* while an immunogold assay may provide the same information, and cellular localization information *in vivo*. It is therefore important to know both the assertion and the evidence to understand what sort of information may be gleaned from the assay. To understand which types of assertions are made in the top-50 high throughput articles, we manually assigned Evidence Codes Ontology (ECO) assertion and evidence terms to the top-50 articles. The ECO ontology is more elaborate than the evidence codes used by Uniprot-GOA, and currently, it is not routinely used in UniprotKB. The results are shown in Table 5.

Interestingly, the third most-frequently used assertion in the top experimental articles was not an experimental assertion, but rather a computational one: the term ECO:00053 “computational combinatorial evidence” is defined as “A type of combinatorial analysis where data are combined and evaluated by an algorithm.” This is not a computational prediction per-se, but rather a combination of several experimental lines of evidence used in a article.

The most used experimental ECO term was ECO:000160 “protein separation fol-

lowed by fragment identification evidence”, which encompasses different types of mass-spectrometry experiments. The next ranking assertion terms were computational: “motif similarity evidence” and “sequence similarity evidence used in automatic assertion”. Those were generally combined with the mass-spectrometry experiments to identify protein sequence fragments reconstructed from the mass-spectrometry. Another frequently used experimental techniques was “RNAi experimental evidence”. This type of experiment was mostly with the articles that used RNA interference in studying *C. elegans*, whose study comprised 12 of the top-50 articles.

Discussion

We have identified several annotation biases in UniProt-GOA. These biases stem from the uneven number of annotations produced by different types of experiments. It is clear that results from high-throughput experiments contribute substantially to the function annotation landscape, as up to 20% of experimentally annotated proteins are annotated by high-throughput assays, with most of them not being annotated by medium- or low-throughput experiments.

At the same time, high throughput experiments produce less information per protein than moderate-, low- and single-throughput experiments as evidenced by the type of terms produced in the Molecular Function and Biological Process ontologies. Furthermore, the number of total GO terms used in the high-throughput experiments is much lower than that used in low and medium throughput experiments. Therefore, while high throughput experiments provide a high coverage of protein function space, it is the low throughput experiments that provide more specific information, as well as a larger diversity of terms.

We have also identified several types of biases that are contributed by high throughput

experiments. First, there is the enrichment of low-information content GO-terms, which means that our understanding of the protein function as provided by high-throughput experiments is more limited than that provided by low-throughput experiments. Second, there is the small number of terms used, when considering the large number of proteins that are being annotated. Third is the general ontology bias towards the cellular component ontology and, to a lesser extent, the Biological Process ontology; at the same time, there are very few articles that deal with the Molecular Function ontology. These biases all stem from the inherent capabilities and limitations of the high-throughput experiments. A fourth, related bias is the organism studied: taken together, studies of *C. elegans* and *A. thaliana* studies comprise 36 of the top-50 annotating articles, or 72%.

The most frequent experiment performed is cell fractionation and mass-spectrometry to assign a Cellular Component ontology terms, and identify the proteins, respectively. Consequently this means that the assignment procedure is limited to the cellular compartments that can be identified with the fractionation methods used. So while Cellular Component is the most frequent annotation used, fractionation and mass-spectrometry is the most common method used to localize proteins in subcellular compartments. A notable exception to the use of fractionation and MS for protein localization is in the top annotating article [13] which uses microscopy for subcellular localization. The only MS experiment in the top-50 articles whose proteins were not annotated with cellular localization was “Proteome survey reveals modularity of the yeast cell machinery” [13]. The resulting annotation was “protein binding” from the Molecular Function ontology. A more detailed discussion on this study follows in the section **Information Capture** below.

The second most frequent type of experiments was RNA Interference (RNAi) whole-genome gene knockdowns in *C. elegans*, *D. melanogaster* and one in *C. albicans*. RNAi

experiments typically use targeted dsRNA which is delivered to the organism and silences specific genes. Typically the experiments here used libraries of RNAi targeted to the whole exome. The phenotypes searched for were mostly associated with embryonic and post-embryonic development. Some studies focused on mitotic spindle assembly [14], lipid storage [14] and endocytic traffic [14]. One study used RNAi to identify mitochondrial protein localization [15]. These studies mostly use the same RNAi libraries, and target the whole *C. elegans* genome using common data resources. Hence the large redundancy observed for *C. elegans* in Table 3.

These two types of assays (mass-spectrometry and RNAi) were strongly linked to the other frequently used experimental ECO terms, by the nature of the methodology used. Thus, “protein separation followed by fragment identification evidence” is usually accompanied with “cell fractionation evidence” and “Western blot evidence”. It should be noted that Western blots are used to verify protein purity rather than specifically to determine protein function. “RNAi experimental evidence” is generally associated with “mutant phenotype evidence used in manual assertion”. All experiments are associated with computational ECO terms, which describe sequence similarity and motif recognition techniques used to identify the sequences found. A strong reliance on computational annotation is therefore an integral part of high throughput experiments. It should be noted that computational annotation here is not necessarily used directly for functional annotation, but rather for identifying the protein by a sequence or motif similarity search.

Information Capture and Scope of GO

We have discussed the information loss that is characteristic of high-throughput experiment, as shown in Figure 4. However, another reason for information loss is the inability to capture certain types of information using the Gene Ontology. GO is purposefully

337 limited to three aspects (MF, BP and CC) of biological function, which are assigned per
 338 protein. However, other aspects of function may emerge from experiments that cannot
 339 be captured by GO. Of note is the study mentioned earlier, “Proteome survey reveals
 340 modularity of the yeast cell machinery” [13]. In this study, the information produced was
 341 primarily of protein complexes, which proteins are binding which proteins, and the rela-
 342 tionship to cellular compartmentalization and biological networks. At the same time, the
 343 only GO-term captured in the curation of this study was “protein binding”. Some, but
 344 not all of this information can be captured more specifically using the children of the term
 345 “protein binding”, but such a process is arguably laborious by manual curation of a high
 346 throughput article. Furthermore, the main information conveyed by this article, namely
 347 the types of protein complexes discovered and how they relate to cellular networks, is out-
 348 side the scope of GO. It is important to realize that while high-throughput experiments
 349 do convey less information per protein within the functional scope as defined by GO,
 350 they still convey composite information such as possible pathway mappings – information
 351 which needs to be captured into annotation databases by means other than GO. In the
 352 example above, the information can be captured by a protein interaction database, but
 353 not by GO annotation.

354 Conclusions

355 Taken together, the annotation biases noted in this study affect our understanding of
 356 protein function space. This, in turn, affects our ability to properly understand the con-
 357 nection between predictors of protein function and the actual function – the hallmark of
 358 computational function annotation. As a dramatic example, during the Critical Assess-
 359 ment of Function Annotation experiment (Radivojac *et al* in press) we have noticed that

roughly 20% of the proteins participating in the challenge and annotated with the Molecular Function Ontology were annotated as “protein binding”, a GO-term that conveys little information. Furthermore, it was shown that the major contribution of “protein binding” term to the CAFA challenge data set was due to high-throughput assays. This illustrates how the concentration of a large number of annotations in a small number of studies provides only a partial picture of the function of these proteins. As we have seen, the picture provided from high throughput experiments is mainly of: 1. subcellular localization cell fractionation and MS based localization and 2. developmental phenotypes. While these data are important, we should be mindful of this bias when examining protein function in the database, even those annotations deemed to be of high quality, i.e. with experimental verification. Furthermore, such a large bias in prior probabilities can adversely affect programs employing prior probabilities, as most machine-learning programs do. Many researchers use programs based on machine-learning algorithms to predict the function of proteins. If the training set for these programs has included a disproportional number of annotations by high-throughput experiments, the results these programs provide will be strongly biased towards a few frequent and shallow GO-terms. In a recent paper, Škunca *et al.* have compared the quality of experimental annotations in UniProtKB, to automated ones, using experimental codes [16]. This study concluded that “The reliability of electronic annotations rivals that of non-experimental curated annotations”. However, that may simply be because of the dominance of high-throughput experiments, with the limited number of GO terms they use, in the experimental annotation landscape.

Several steps can be taken to remedy this situation. Annotations are derived from high-throughput experiments can be flagged as such in the database. The flagging can then be read by sequence similarity or other search software, and flagged proteins removed or otherwise marked in the search. In a typical scenario, a researcher will BLAST their query

385 protein to determine its function by sequence similarity. If a target protein is tagged as
 386 annotated by a high throughput assay, it would be removed from the search if asked to do
 387 so by the user. This filtering can also be done by assay type, number of proteins annotated
 388 per experiment, or a combination of the above. This requires that GO-annotated proteins
 389 should also be annotated with assertion codes in addition to the evidence codes and GO
 390 term-codes; but given the large volume of data in UniprotKB is it hard to expect such
 391 massive reannotation with assertion terms undertaken. (Any other ideas?)

392 We call upon the communities of annotators, computational biologists and experi-
 393 mental biologists to be mindful of the phenomenon of the high-throughput experimental
 394 biases described in this study, and to work to understand its implications and mitigate
 395 its impact.

396 **Note on Methods**

397 We used the UniProtKB-GOA database from December 2011. Data analyses were per-
 398 formed using Python scripts. ECO terms classifying the proteins in the top 50 experiments
 399 were assigned to the proteins manually after reading the articles. All data and scripts are
 400 available on: <http://github.com/FriedbergLab/DataBias/>

401 **Acknowledgments**

402 We thank Predrag Radivojac, Nives Škunca, Cristophe Dessimoz and members of the
 403 Friedberg and Babbitt labs for insightful discussions.

404 Funding

405 This research was funded, in part by NSF / ABI XXXXXXXX-XXX award to IF and
406 XXXXXXXXXXXXX to PCB.

407 References

408 References

- 409 1. Friedberg I (2006) Automated protein function prediction—the genomic challenge.
410 Brief Bioinform 7: 225–242.
- 411 2. Schnoes AM, Brown SD, Dodevski I, Babbitt PC (2009) Annotation error in public
412 databases: Misannotation of molecular function in enzyme superfamilies. PLoS
413 Comput Biol 5: e1000605+.
- 414 3. Erdin S, Lisewski AM, Lichtarge O (2011) Protein function prediction: towards
415 integration of similarity metrics. Current Opinion in Structural Biology 21: 180 -
416 188.
- 417 4. Rentzsch R, Orengo CA (2009) Protein function prediction the power of multiplic-
418 ity. Trends in Biotechnology 27: 210 - 219.
- 419 5. Sthl PL, Lundberg J (2012) Toward the single-hour high-quality genome. Annual
420 Review of Biochemistry 81: 359-378.
- 421 6. Sboner A, Mu X, Greenbaum D, Auerbach R, Gerstein M (2011) The real cost of
422 sequencing: higher than you think! Genome Biology 12: 125+.

- 423 7. Dimmer EC, Huntley RP, Alam-Faruque Y, Sawford T, O'Donovan C, et al. (2012)
424 The uniprot-go annotation database in 2011. *Nucleic Acids Research* 40: D565–
425 D570.
- 426 8. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, et al. (2000) Gene ontology:
427 tool for the unification of biology. *Nature Genetics* 25: 25–29.
- 428 9. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, et al. (1997) Gapped
429 blast and psi-blast: a new generation of protein database search programs. *Nucleic
430 acids research* 25: 3389–3402.
- 431 10. Li W, Godzik A (2006) Cd-hit: a fast program for clustering and comparing large
432 sets of protein or nucleotide sequences. *Bioinformatics* 22: 1658–1659.
- 433 11. Lord PW, Stevens RD, Brass A, Goble CA (2003) Investigating semantic simi-
434 larity measures across the gene ontology: the relationship between sequence and
435 annotation. *Bioinformatics* 19: 1275–1283.
- 436 12. Pesquita C, Faria D, Falcão AO, Lord P, Couto FM (2009) Semantic similarity in
437 biomedical ontologies. *PLoS Comput Biol* 5: e1000443+.
- 438 13. Barbe L, Lundberg E, Oksvold P, Stenius A, Lewin E, et al. (2008) Toward a
439 confocal subcellular atlas of the human proteome. *Mol Cell Proteomics* 7: 499–
440 508.
- 441 14. Goshima G, Wollman R, Goodwin SS, Zhang N, Scholey JM, et al. (2007) Genes
442 required for mitotic spindle assembly in *Drosophila* S2 cells. *Science* 316: 417–421.

- 443 15. Hughes JR, Meireles AM, Fisher KH, Garcia A, Antrobus PR, et al. (2008) A
444 microtubule interactome: complexes with roles in cell cycle and mitosis. PLoS Biol
445 6: e98.
- 446 16. Škunca N, Altenhoff A, Dessimoz C (2012) Quality of computationally inferred
447 gene ontology annotations. PLoS Comput Biol 8: e1002533+.
- 448 17. Matsuyama A, Arai R, Yashiroda Y, Shirai A, Kamata A, et al. (2006) ORFeome
449 cloning and global analysis of protein localization in the fission yeast *Schizosaccha-*
450 *romyces pombe*. Nat Biotechnol 24: 841–847.
- 451 18. Pagliarini DJ, Calvo SE, Chang B, Sheth SA, Vafai SB, et al. (2008) A mitochon-
452 drial protein compendium elucidates complex I disease biology. Cell 134: 112–123.
- 453 19. Simmer F, Moorman C, van der Linden AM, Kuijk E, van den Berghe PV, et al.
454 (2003) Genome-wide RNAi of *C. elegans* using the hypersensitive *rrf-3* strain reveals
455 novel gene functions. PLoS Biol 1: E12.
- 456 20. Huh WK, Falvo JV, Gerke LC, Carroll AS, Howson RW, et al. (2003) Global
457 analysis of protein localization in budding yeast. Nature 425: 686–691.
- 458 21. Zybaylov B, Rutschow H, Friso G, Rudella A, Emanuelsson O, et al. (2008) Sorting
459 signals, N-terminal modifications and abundance of the chloroplast proteome. PLoS
460 ONE 3: e1994.
- 461 22. Sonnichsen B, Koski LB, Walsh A, Marschall P, Neumann B, et al. (2005) Full-
462 genome RNAi profiling of early embryogenesis in *Caenorhabditis elegans*. Nature
463 434: 462–469.

- 464 23. Mootha VK, Bunkenborg J, Olsen JV, Hjerrild M, Wisniewski JR, et al. (2003)
465 Integrated analysis of protein composition, tissue diversity, and gene regulation in
466 mouse mitochondria. *Cell* 115: 629–640.
- 467 24. Benschop JJ, Mohammed S, O’Flaherty M, Heck AJ, Slijper M, et al. (2007) Quan-
468 titative phosphoproteomics of early elicitor signaling in Arabidopsis. *Mol Cell Pro-*
469 *teomics* 6: 1198–1214.
- 470 25. Kamath RS, Fraser AG, Dong Y, Poulin G, Durbin R, et al. (2003) Systematic
471 functional analysis of the *Caenorhabditis elegans* genome using RNAi. *Nature* 421:
472 231–237.
- 473 26. Mawuenyega KG, Forst CV, Dobos KM, Belisle JT, Chen J, et al. (2005) My-
474 cobacterium tuberculosis functional network analysis by global subcellular protein
475 profiling. *Mol Biol Cell* 16: 396–404.
- 476 27. Ito J, Batth TS, Petzold CJ, Redding-Johanson AM, Mukhopadhyay A, et al. (2011)
477 Analysis of the Arabidopsis cytosolic proteome highlights subcellular partitioning
478 of central plant metabolism. *J Proteome Res* 10: 1571–1582.
- 479 28. Rual JF, Ceron J, Koreth J, Hao T, Nicot AS, et al. (2004) Toward improving
480 *Caenorhabditis elegans* phenome mapping with an ORFeome-based RNAi library.
481 *Genome Res* 14: 2162–2168.
- 482 29. Reinders J, Zahedi RP, Pfanner N, Meisinger C, Sickmann A (2006) Toward the
483 complete yeast mitochondrial proteome: multidimensional separation techniques
484 for mitochondrial proteomics. *J Proteome Res* 5: 1543–1554.
- 485 30. Fernandez-Calvino L, Faulkner C, Walshaw J, Saalbach G, Bayer E, et al. (2011)
486 Arabidopsis plasmodesmal proteome. *PLoS ONE* 6: e18880.

- 487 31. Gu S, Chen J, Dobos KM, Bradbury EM, Belisle JT, et al. (2003) Comprehensive
488 proteomic profiling of the membrane constituents of a *Mycobacterium tuberculosis*
489 strain. *Mol Cell Proteomics* 2: 1284–1296.
- 490 32. Ferro M, Brugiére S, Salvi D, Seigneurin-Berny D, Court M, et al. (2010)
491 ATCHLORO, a comprehensive chloroplast proteome database with subplastidial
492 localization and curated information on envelope proteins. *Mol Cell Proteomics* 9:
493 1063–1084.
- 494 33. Kleffmann T, Russenberger D, von Zychlinski A, Christopher W, Sjolander K, et al.
495 (2004) The *Arabidopsis thaliana* chloroplast proteome reveals pathway abundance
496 and novel protein functions. *Curr Biol* 14: 354–362.
- 497 34. Sassetti CM, Boyd DH, Rubin EJ (2003) Genes required for mycobacterial growth
498 defined by high density mutagenesis. *Mol Microbiol* 48: 77–84.
- 499 35. Balklava Z, Pant S, Fares H, Grant BD (2007) Genome-wide analysis identifies
500 a general requirement for polarity proteins in endocytic traffic. *Nat Cell Biol* 9:
501 1066–1073.
- 502 36. Mitra SK, Gantt JA, Ruby JF, Clouse SD, Goshe MB (2007) Membrane proteomic
503 analysis of *Arabidopsis thaliana* using alternative solubilization techniques. *J Pro-*
504 *teome Res* 6: 1933–1950.
- 505 37. Maeda I, Kohara Y, Yamamoto M, Sugimoto A (2001) Large-scale analysis of gene
506 function in *Caenorhabditis elegans* by high-throughput RNAi. *Curr Biol* 11: 171–
507 176.
- 508 38. Ceron J, Rual JF, Chandra A, Dupuy D, Vidal M, et al. (2007) Large-scale RNAi
509 screens identify novel genes that interact with the *C. elegans* retinoblastoma path-

- 510 way as well as splicing-related components with synMuv B activity. BMC Dev Biol
511 7: 30.
- 512 39. Sickmann A, Reinders J, Wagner Y, Joppich C, Zahedi R, et al. (2003) The pro-
513 teome of *Saccharomyces cerevisiae* mitochondria. Proc Natl Acad Sci USA 100:
514 13207–13212.
- 515 40. Gavin AC, Aloy P, Grandi P, Krause R, Boesche M, et al. (2006) Proteome survey
516 reveals modularity of the yeast cell machinery. Nature 440: 631–636.
- 517 41. Green RA, Kao HL, Audhya A, Arur S, Mayers JR, et al. (2011) A high-resolution
518 *C. elegans* essential gene network based on phenotypic profiling of a complex tissue.
519 Cell 145: 470–482.
- 520 42. Simpson JC, Wellenreuther R, Poustka A, Pepperkok R, Wiemann S (2000) Sys-
521 tematic subcellular localization of novel proteins identified by large-scale cDNA
522 sequencing. EMBO Rep 1: 287–292.
- 523 43. Marmagne A, Ferro M, Meinnel T, Bruley C, Kuhn L, et al. (2007) A high content
524 in lipid-modified peripheral proteins and integral receptor kinases features in the
525 arabidopsis plasma membrane proteome. Mol Cell Proteomics 6: 1980–1996.
- 526 44. Dunkley TP, Hester S, Shadforth IP, Runions J, Weimar T, et al. (2006) Mapping
527 the Arabidopsis organelle proteome. Proc Natl Acad Sci USA 103: 6518–6523.
- 528 45. Jaquinod M, Villiers F, Kieffer-Jaquinod S, Hugouvieux V, Bruley C, et al. (2007)
529 A proteomics dissection of Arabidopsis thaliana vacuoles isolated from cell culture.
530 Mol Cell Proteomics 6: 394–412.

- 531 46. Heazlewood JL, Tonti-Filippini JS, Gout AM, Day DA, Whelan J, et al. (2004) Ex-
532 perimental analysis of the Arabidopsis mitochondrial proteome highlights signaling
533 and regulatory components, provides assessment of targeting prediction programs,
534 and indicates plant-specific mitochondrial proteins. *Plant Cell* 16: 241–256.
- 535 47. Ashrafi K, Chang FY, Watts JL, Fraser AG, Kamath RS, et al. (2003) Genome-
536 wide RNAi analysis of *Caenorhabditis elegans* fat regulatory genes. *Nature* 421:
537 268–272.
- 538 48. Piano F, Schetter AJ, Morton DG, Gunsalus KC, Reinke V, et al. (2002) Gene
539 clustering based on RNAi phenotypes of ovary-enriched genes in *C. elegans*. *Curr*
540 *Biol* 12: 1959–1964.
- 541 49. Carter C, Pan S, Zouhar J, Avila EL, Girke T, et al. (2004) The vegetative vacuole
542 proteome of *Arabidopsis thaliana* reveals predicted and unexpected proteins. *Plant*
543 *Cell* 16: 3285–3303.
- 544 50. Da Cruz S, Xenarios I, Langridge J, Vilbois F, Parone PA, et al. (2003) Proteomic
545 analysis of the mouse liver mitochondrial inner membrane. *J Biol Chem* 278: 41566–
546 41571.
- 547 51. Rual JF, Venkatesan K, Hao T, Hirozane-Kishikawa T, Dricot A, et al. (2005)
548 Towards a proteome-scale map of the human protein-protein interaction network.
549 *Nature* 437: 1173–1178.
- 550 52. Bakthavatsalam D, Gomer RH (2010) The secreted proteome profile of developing
551 *Dictyostelium discoideum* cells. *Proteomics* 10: 2556–2559.
- 552 53. Froehlich JE, Wilkerson CG, Ray WK, McAndrew RS, Osteryoung KW, et al.
553 (2003) Proteomic study of the *Arabidopsis thaliana* chloroplastic envelope mem-

- 554 brane utilizing alternatives to traditional two-dimensional electrophoresis. *J Pro-*
555 *teome Res* 2: 413–425.
- 556 54. Stroschein-Stevenson SL, Foley E, O’Farrell PH, Johnson AD (2006) Identification
557 of *Drosophila* gene products required for phagocytosis of *Candida albicans*. *PLoS*
558 *Biol* 4: e4.
- 559 55. Rutschow H, Ytterberg AJ, Friso G, Nilsson R, van Wijk KJ (2008) Quantitative
560 proteomics of a chloroplast SRP54 sorting mutant and its genetic interactions with
561 CLPC1 in *Arabidopsis*. *Plant Physiol* 148: 156–175.
- 562 56. Kumar A, Agarwal S, Heyman JA, Matson S, Heidtman M, et al. (2002) Subcellular
563 localization of the yeast proteome. *Genes Dev* 16: 707–719.
- 564 57. Fraser AG, Kamath RS, Zipperlen P, Martinez-Campos M, Sohrmann M, et al.
565 (2000) Functional genomic analysis of *C. elegans* chromosome I by systematic RNA
566 interference. *Nature* 408: 325–330.
- 567 58. Gonczy P, Echeverri C, Oegema K, Coulson A, Jones SJ, et al. (2000) Functional
568 genomic analysis of cell division in *C. elegans* using RNAi of genes on chromosome
569 III. *Nature* 408: 331–336.
- 570 59. Suzuki H, Fukunishi Y, Kagawa I, Saito R, Oda H, et al. (2001) Protein-protein
571 interaction panel using mouse full-length cDNAs. *Genome Res* 11: 1758–1765.
- 572 60. Sarry JE, Kuhn L, Ducruix C, Lafaye A, Junot C, et al. (2006) The early re-
573 sponses of *Arabidopsis thaliana* cells to cadmium exposure explored by protein and
574 metabolite profiling analyses. *Proteomics* 6: 2180–2198.

- 575 61. Chen D, Toone WM, Mata J, Lyne R, Burns G, et al. (2003) Global transcriptional
576 responses of fission yeast to environmental stress. *Mol Biol Cell* 14: 214–229.
- 577 62. Herold N, Will CL, Wolf E, Kastner B, Urlaub H, et al. (2009) Conservation of the
578 protein composition and electron microscopy structure of *Drosophila melanogaster*
579 and human spliceosomal complexes. *Mol Cell Biol* 29: 281–301.
- 580 63. Bayer EM, Bottrill AR, Walshaw J, Vigouroux M, Naldrett MJ, et al. (2006) Ara-
581 bidopsis cell wall proteome defined using multidimensional protein identification
582 technology. *Proteomics* 6: 301–311.

583 Figures

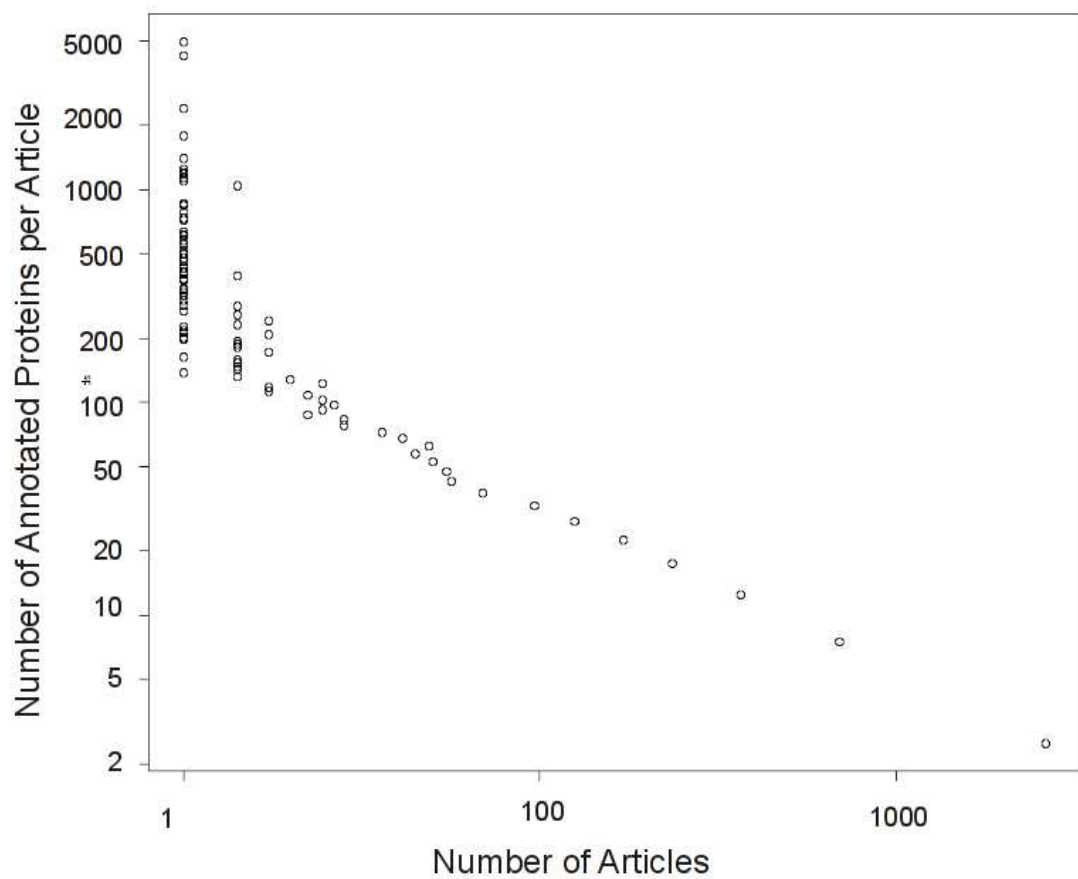


Figure 1. Distribution of the number of proteins annotated per article.

X-axis: number of annotating articles. Y-axis: number of annotated proteins. The distribution was found to be logarithmic with a significant ($R^2 = 0.72$; $p < 1.10 \times 10^{-18}$) linear fit to the log-log plot. The data came from 76137 articles annotating 256033 proteins with GO experimental evidence codes, in Uniprot-GOA 12/2011.

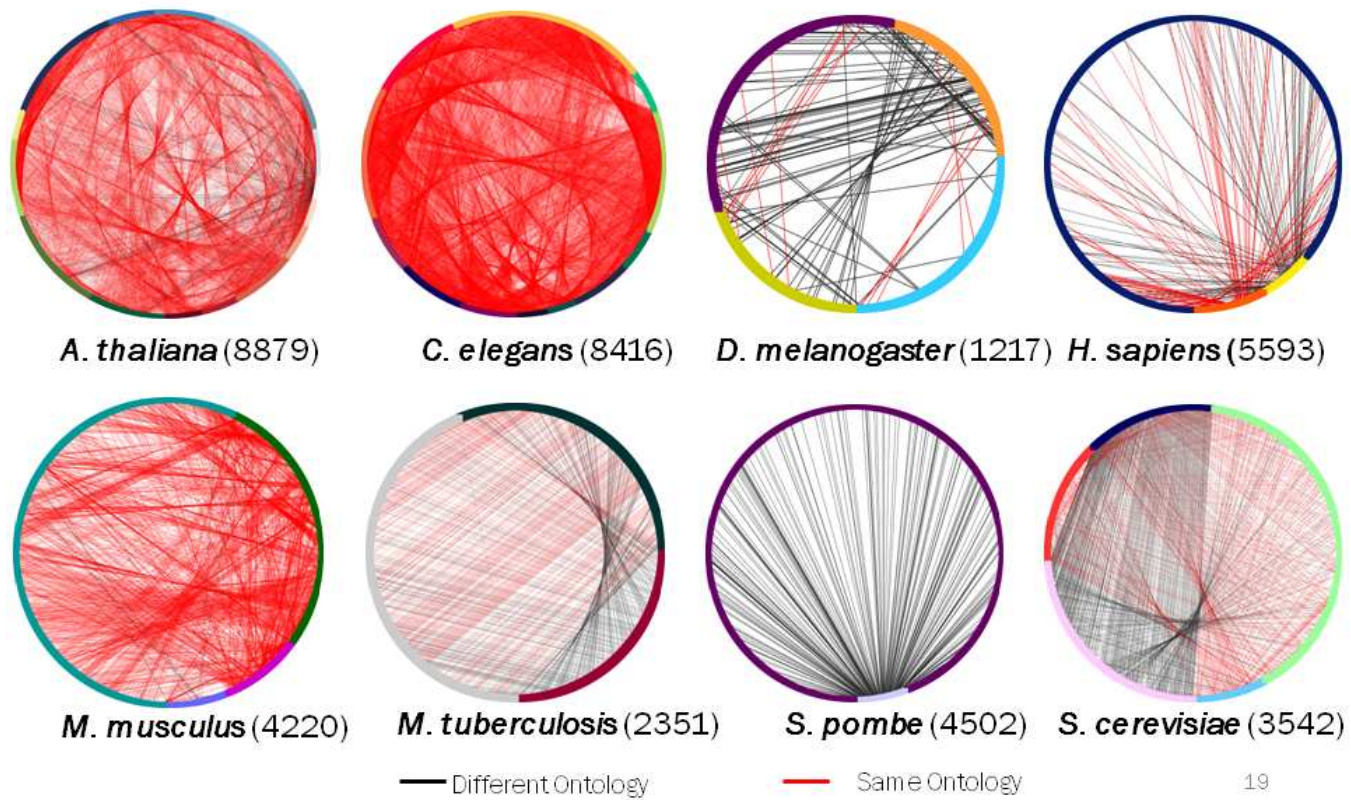


Figure 2. Redundancy in proteins described by the top-50 articles. A circle represents the sum total of articles annotating each organism. Each colored arch is composed of all the proteins in a single article. A line is drawn between any two points on the circle if the proteins they represent have 100% sequence identity. A black line is drawn if they are annotated with a different ontology (e.g. in one article the protein is annotated with the MFO, and in another article with BPO); a red line if they are annotated in the same ontology. Example: *S. pombe* is described by two articles, one with few protein (light arch on bottom) and one with many (dark arch encompassing most of circle). Many of the same proteins are annotated by both articles. See table 3 for numbers.

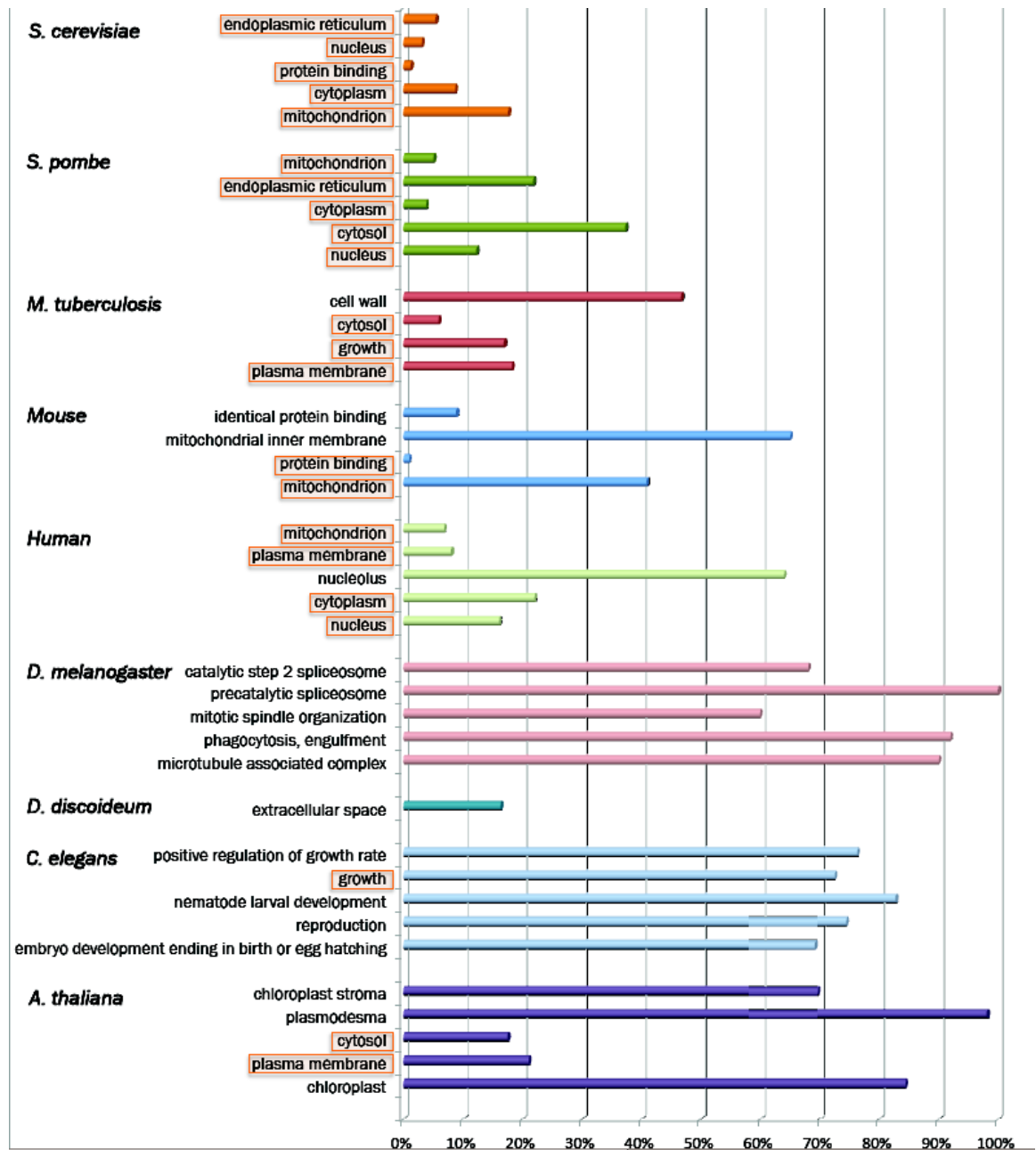


Figure 3. Relative contribution of top-50 articles to the annotation of major model organisms. The length of each bar represents the percentage of proteins annotated by the top-50 articles in a given organism by a given GO term.

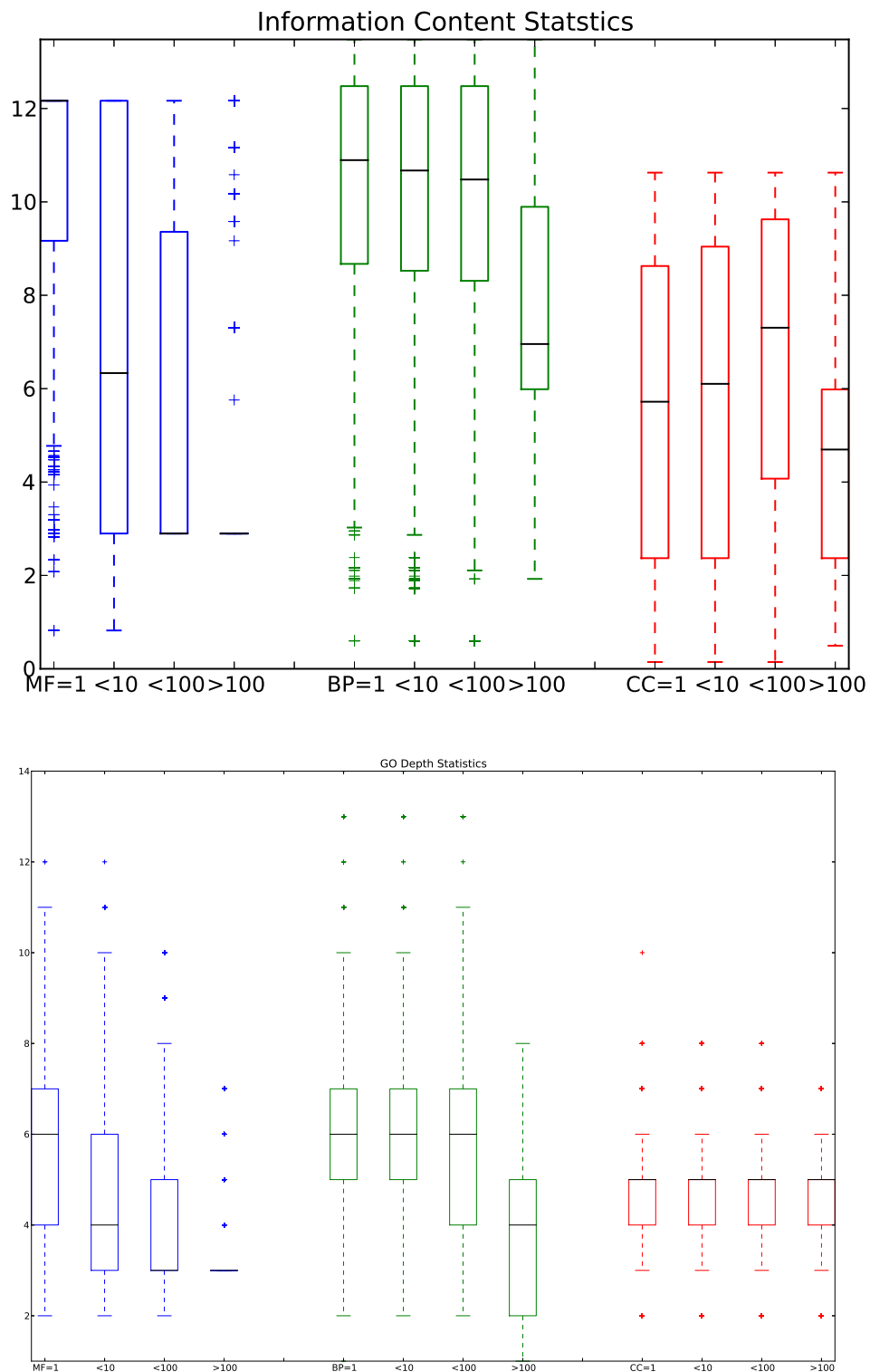


Figure 4. Information provided by articles depending on the number of proteins the articles annotate. Articles are grouped into cohorts: 1: one protein annotated by article; <10 : more than 1, less than 10 annotated; <100 : more than 10, less than 100 annotated; ≥ 100 : more than 100 proteins annotated per article. Blue bars: Molecular Function ontology; Green bars: Biological Process ontology; Red bars: Cellular Component ontology. Information is gauged by **A**: Information Content and **B**: GO depth. See text for details.

Table 1. Top 50 Annotating Articles

N	Proteins	Annotations	Species	ref.	MFO	BPO	CCO
1	4937	11050	<i>H. sapiens</i>	[13]	0	0	11050
2	4247	7046	<i>S. pombe</i>	[17]	0	0	7046
3	2412	2412	<i>H. sapiens</i>	[18]	0	0	2412
4	1791	5918	<i>C. elegans</i>	[19]	0	5918	0
5	1406	1863	<i>S. cerevisiae</i>	[20]	0	0	1863
6	1251	1251	<i>A. thaliana</i>	[21]	0	0	1251
7	1205	1476	<i>C. elegans</i>	[22]	0	1476	0
8	1186	1213	<i>M. musculus</i>	[23]	0	0	1213
9	1136	1136	<i>A. thaliana</i>	[24]	0	0	1136
10	1101	2269	<i>C. elegans</i>	[25]	0	2269	0
11	1043	1365	<i>M. tuberculosis</i>	[26]	0	0	1365
12	1041	1041	<i>A. thaliana</i>	[27]	0	0	1041
13	865	1533	<i>C. elegans</i>	[28]	0	1533	0
14	845	845	<i>S. cerevisiae</i>	[29]	0	0	845
15	784	784	<i>A. thaliana</i>	[30]	0	0	784
16	735	735	<i>M. tuberculosis</i>	[31]	0	0	735
17	724	882	<i>A. thaliana</i>	[32]	0	0	882
18	634	634	<i>A. thaliana</i>	[33]	0	0	634
19	613	613	Mycobacter sp.	[34]	0	613	0
20	607	661	<i>C. elegans</i>	[35]	0	659	2

Continued on next page

N	Proteins	Annotations	Species	ref.	MFO	BPO	CCO
21	577	577	<i>A. thaliana</i>	[36]	0	0	577
22	553	884	<i>C. elegans</i>	[37]	0	884	0
23	516	5972	<i>C. elegans</i>	[38]	0	5972	0
24	503	503	<i>S. cerevisiae</i>	[39]	0	0	503
25	498	638	<i>S. cerevisiae</i>	[40]	638	0	0
26	479	848	<i>C. elegans</i>	[41]	0	848	0
27	465	468	<i>H. sapiens</i>	[42]	0	0	468
28	436	436	<i>A. thaliana</i>	[43]	0	0	436
29	430	513	<i>A. thaliana</i>	[44]	0	0	513
30	413	456	<i>D. melanogaster</i>	[15]	0	39	417
31	401	401	<i>A. thaliana</i>	[45]	0	0	401
32	392	392	<i>A. thaliana</i>	[46]	0	0	392
33	392	639	<i>C. elegans</i>	[47]	0	639	0
34	383	917	<i>C. elegans</i>	[48]	0	917	0
35	380	380	<i>A. thaliana</i>	[49]	0	0	380
36	375	375	<i>M. musculus</i>	[50]	0	0	375
37	343	509	<i>H. sapiens</i>	[51]	509	0	0
38	338	338	Ddiscoideum	[52]	0	0	338
39	328	328	<i>A. thaliana</i>	[53]	0	0	328
40	319	329	<i>C. albicans</i>	[54]	1	328	0
41	305	312	<i>A. thaliana</i>	[55]	0	0	312
42	290	331	<i>S. cerevisiae</i>	[56]	0	0	331

Continued on next page

N	Proteins	Annotations	Species	ref.	MFO	BPO	CCO
43	285	761	<i>C. elegans</i>	[57]	0	761	0
44	283	499	<i>C. elegans</i>	[58]	0	499	0
45	266	433	<i>M. musculus</i>	[59]	433	0	0
46	260	260	<i>A. thaliana</i>	[60]	0	260	0
47	258	259	<i>S. pombe</i>	[61]	0	259	0
48	244	397	<i>D. melanogaster</i>	[14]	0	367	30
49	242	397	<i>D. melanogaster</i>	[62]	0	0	397
50	241	263	<i>A. thaliana</i>	[63]	0	0	263

585 **The top 50 annotating articles.** **N**: article rank; **Proteins**: number of proteins
 586 annotated in this article; **Annotations**: number of annotating GO terms; **Species**:
 587 annotated species; **ref.** annotating article; **MFO/BPO/CCO**: number of proteins
 588 annotated in the Molecular Function, Biological Process and Cellular Component
 589 ontologies, respectively.

Table 2. Annotation Cohorts

Articles annotating the following number of proteins	1	$1 < n \leq 10$	$10 < n \leq 100$	$n > 100$	SUM
Number of proteins annotated	20699	46383	26485	31411	124978
Number of annotating articles	41156	32201	2672	108	76137
Percent of proteins annotated	16.56	37.11	21.19	25.13	100
Percent of annotating articles	54.09	42.32	3.51	0.14	100

Table caption

Table 3. Sequence Redundancy in Top-50 Annotating Articles

Species	num. articles	num. prot	Clusters at 100%	% redundancy	Mean genes/ cluster
<i>C. elegans</i>	12	8416	3338	60	3.74
<i>A. thaliana</i>	16	8879	4694	47	3.92
<i>M. musculus</i>	3	4220	2273	46	2.75
<i>M. tuberculosis</i>	2	2351	1702	28	2.22
<i>S. cerevisiae</i>	5	3542	2550	28	2.33
<i>H. sapiens</i>	4	5593	4509	19	2.36
<i>D. melanogaster</i>	3	1217	1003	18	2.17
<i>S. pombe</i>	2	4502	4281	5	2.00

Species: annotated species; **num. articles** number of annotating articles; **num. prot:** number of proteins annotated by top-50 articles for that species; **Clusters at 100%:** number of clusters of 100% identical proteins; **% redundancy:** the ratio between column 3 and column 2: this is the percentage of proteins annotated more than once for a given species in the top 50 articles; **Mean genes/cluster:** the mean number of genes per cluster, for clusters having more than a single gene.

Table 4. Annotation Consistency in Top 50 articles

Species	Ont.	num prot	mean $k_{P,O}$	stdv	stderr	num articles	num terms
A. thaliana	CCO	1941	0.251	0.328	0.007	15	18
C. elegans	BPO	1847	0.388	0.239	0.006	12	41
D. melanogaster	BPO	76	0.086	0.22	0.025	3	8
D. melanogaster	CCO	81	0.068	0.234	0.026	3	5
H. sapiens	CCO	167	0.285	0.365	0.028	2	20
M. musculus	CCO	807	0.832	0.291	0.01	3	2
S. cerevisiae	CCO	744	0.759	0.379	0.014	4	15
B. tuberculosis	CCO	532	0.309	0.41	0.018	2	3

Species: annotated species; **Ontology:** annotating GO ontology; **num prot:** number of annotated proteins in that species & ontology that are annotated by more than one paper. **mean, stdv, stderr:** mean number of consistent annotations for a protein in that species and ontology, standard deviation from the mean and standard error. **num articles:** number of annotating articles **num terms** number of annotating terms. Annotations by less than two articles or two terms (or both) for the same protein/ontology combination have been omitted.

Table 6. Assertion codes used in top-50 papers

PMID	Ref.	ECO ID's
14551910	[19]	ECO:0000019 ECO:0000315
15791247	[22]	ECO:0000019 ECO:0000315
12529643	[47]	ECO:0000019 ECO:0000315
20061580	[32]	ECO:0000160 ECO:0000004 ECO:0000250 ECO:0000028 ECO:0000081 ECO:0000083 ECO:0000112
12529635	[25]	ECO:0000019 ECO:0000315 ECO:0000031 ECO:0000028
18433294	[15]	ECO:0000160 ECO:0000019 ECO:0000004 ECO:0000112 ECO:0000249 ECO:0000315 ECO:0000053
14651853	[23]	ECO:0000160 ECO:0000126 ECO:0000010 ECO:0000245 ECO:0000287
11914276	[56]	ECO:0000015 ECO:0000092 ECO:0000007 ECO:0000028 ECO:0000083 ECO:0000053
21529718	[41]	ECO:0000315 ECO:0000019 ECO:0000053
12657046	[34]	ECO:0000315 ECO:0000015 ECO:0000097 ECO:0000053
18431481	[21]	ECO:0000160 ECO:0000126 ECO:0000081 ECO:0000004
14562095	[20]	ECO:0000126 ECO:0000124

Continued on next page

PMID	Ref.	ECO ID's
12445391	[48]	ECO:0000019 ECO:0000315 ECO:0000250 ECO:0000053
18981222	[62]	ECO:0000208 ECO:0000160 ECO:0000249 ECO:0000181
11256614	[42]	ECO:0000053 ECO:0000249 ECO:0000028 ECO:0000126 ECO:0000128
11099033	[57]	ECO:0000019 ECO:0000315 ECO:0000053 ECO:0000031 ECO:0000245
17412918	[14]	ECO:0000019 ECO:0000315
16336044	[54]	ECO:0000019 ECO:0000315
16502469	[60]	ECO:0000160 ECO:0000249 ECO:0000106 ECO:0000108
12529438	[61]	ECO:0000104 ECO:0000266
11099034	[58]	ECO:0000019 ECO:0000315 ECO:0000031
15489339	[28]	ECO:0000019 ECO:0000315 ECO:0000176
21166475	[27]	ECO:0000160 ECO:0000004 ECO:0000053 ECO:0000031 ECO:0000112
15525680	[26]	ECO:0000160 ECO:0000004 ECO:0000053
16618929	[44]	ECO:0000249 ECO:0000028 ECO:0000160 ECO:0000004
18029348	[13]	ECO:0000324 ECO:0000092 ECO:0000007

Continued on next page

PMID	Ref.	ECO ID's
17151019	[45]	ECO:0000160 ECO:0000004 ECO:0000112 ECO:0000249
21533090	[30]	ECO:0000160 ECO:0000004 ECO:0000249 ECO:0000028
17644812	[43]	ECO:0000053 ECO:0000249 ECO:0000028 ECO:0000081 ECO:0000004
17432890	[36]	ECO:0000053 ECO:0000160 ECO:0000004 ECO:0000249 ECO:0000081 ECO:0000044
17317660	[24]	ECO:0000160 ECO:0000004 ECO:0000245 ECO:0000249
15028209	[33]	ECO:0000160 ECO:0000004 ECO:0000028 ECO:0000053 ECO:0000287 ECO:0000044 ECO:0000250 ECO:0000081
17704769	[35]	ECO:0000315 ECO:0000019
14532352	[31]	ECO:0000160 ECO:0000004 ECO:0000249 ECO:0000028 ECO:0000083
12938931	[53]	ECO:0000160 ECO:0000249 ECO:0000112 ECO:0000004
16823372	[17]	ECO:0000128 ECO:0000112 ECO:0000122 ECO:0000231
18633119	[55]	ECO:0000160 ECO:0000249 ECO:0000112 ECO:0000004 ECO:0000315

Continued on next page

PMID	Ref.	ECO ID's
17417969	[38]	ECO:0000019 ECO:0000315
18614015	[18]	ECO:0000160 ECO:0000126 ECO:0000004
16189514	[51]	ECO:0000068 ECO:0000053 ECO:0000022
20422638	[52]	ECO:0000160 ECO:0000044 ECO:0000028
15539469	[49]	ECO:0000160 ECO:0000004 ECO:0000249 ECO:0000028
16823961	[29]	ECO:0000160 ECO:0000004 ECO:0000249
16429126	[40]	ECO:0000079 ECO:0000053 ECO:0000160
16287169	[63]	ECO:0000160 ECO:0000249 ECO:0000081 ECO:0000028 ECO:0000083 ECO:0000053 ECO:0000248
14576278	[39]	ECO:0000053 ECO:0000160 ECO:0000004 ECO:0000249 ECO:0000028 ECO:0000083 ECO:0000112
12865426	[50]	ECO:0000160 ECO:0000004 ECO:0000250 ECO:0000028
11591653	[59]	ECO:0000025 ECO:0000112
11231151	[37]	ECO:0000019 ECO:0000315
14671022	[46]	ECO:0000160 ECO:0000004 ECO:0000249 ECO:0000081 ECO:0000044 ECO:0000053

Table 5. Frequency of assertion codes used in top-50 papers

N	ECO id	ECO term	Articles
1	ECO:0000160	protein separation followed by fragment identification evidence	25
2	ECO:0000004	cell fractionation evidence	21
3	ECO:0000053	computational combinatorial evidence	18
4	ECO:0000249	sequence similarity evidence used in automatic assertion	18
5	ECO:0000315	mutant phenotype evidence used in manual assertion	16
6	ECO:0000019	RNAi experimental evidence	15
7	ECO:0000028	motif similarity evidence	14
8	ECO:0000112	Western blot evidence	9
9	ECO:0000081	targeting sequence prediction evidence	7
10	ECO:0000083	transmembrane domain prediction evidence	5
11	ECO:0000126	GFP fusion protein localization evidence	5
12	ECO:0000250	sequence similarity evidence used in manual assertion	4
13	ECO:0000031	protein BLAST evidence used in manual assertion	4
14	ECO:0000044	sequence similarity evidence	4
15	ECO:0000104	microarray RNA expression level evidence	3
16	ECO:0000245	computational combinatorial evidence used in manual assertion	3
17	ECO:0000015	transposon integration	2
18	ECO:0000128	YFP fusion protein localization evidence	2
19	ECO:0000092	epitope-tagged protein immunolocalization evidence	2
20	ECO:0000007	immunofluorescence evidence	2
21	ECO:0000248	sequence alignment evidence used in automatic assertion	1
22	ECO:0000010	protein expression evidence	1
23	ECO:0000231	qRT-PCR evidence	1
24	ECO:0000122	protein localization evidence	1
25	ECO:0000181	<i>in-vitro</i> assay evidence	1
26	ECO:0000208	protein BLAST evidence	1
27	ECO:0000108	reverse transcription polymerase chain reaction transcription evidence	1
28	ECO:0000062	genomic microarray evidence	1
29	ECO:0000106	Northern assay evidence	1
30	ECO:0000026	nucleic acid hybridization evidence	1
31	ECO:0000068	yeast 2-hybrid evidence	1
32	ECO:0000176	mutant visible phenotype evidence	1
33	ECO:0000324	imaging assay evidence	1
34	ECO:0000079	affinity chromatography evidence	1
35	ECO:0000022	co-purification evidence	1
36	ECO:0000266	sequence orthology evidence used in manual assertion	1
37	ECO:0000025	hybrid interaction evidence	1
38	ECO:0000124	RFP fusion protein localization	1

Assertion codes we assigned to the top-50 annotating papers. The table entries are ranked by the frequency of the assignments, i.e. 25 papers are assigned with term ECO:0000160, 21 were assigned ECO:0000004, etc. Entries in **boldface** are for