

1 Biases in the Experimental Annotations of Protein 2 Function and their Effect on Our Understanding of 3 Protein Function Space

4 1 Alexandra Schnoes Dept/Program/Center, University of California, San
5 Francisco, San Francisco, CA, USA

6 2 David C. Ream Department of Microbiology, Miami University, Oxford,
7 OH, USA

8 3 Alexander W. Thorman, Department of Microbiology, Miami University,
9 Oxford, OH, USA

10 4 Patricia C. Babbitt Dept/Program/Center, University of California, San
11 Francisco, San Francisco, CA, USA

12 5 Iddo Friedberg, Department of Microbiology and Computer Science &
13 Software engineering, Miami University, Oxford, OH, USA

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

15 Abstract

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

Should probably not have that
comment in bckground.

24 enabled by proteins.

25 **Results:** We examine the annotation of UniProtKB by the Gene Ontology Annotation
 26 project (GOA), and show that the distribution of proteins per article is a **log-odd**, with
 27 0.06% of articles dominating 20% of the annotations. Since each of the dominant articles
 28 describes the use of an assay that can find only one function or a small group of functions,
 29 this leads to substantial biases, in several aspects, in what we know about the function of
 30 many proteins.

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

36 Author Summary

37 Introduction

38 Functional annotation of proteins is a primary challenge in molecular biology today
 39 [?, 1–3]. The ongoing improvements in sequencing technology had the emphasis shift-
 40 ing from realizing the \$1000 genome to the 1-hour genome [?]. The ability to rapidly and
 41 cheaply sequence genomes is creating a flood of sequence data, which require extensive
 42 analysis and characterization before these data can be useful. A large proportion of this
 43 work involves assigning biological function to newly determined gene sequences, a process
 44 that is both complex and costly [4]. Furthermore, the ability to accurately assign function

Pmid in very old version of paper.

through computational means is challenging and open problem [5]. To aid current annotation procedures and improve computational 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 [6], 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 [7] 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 specific, experimentally studied sequences. 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 [5].

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 few experiments dominate the protein function landscape. In this work we analyzed the relative contribution of articles to the experimental annotations in UniProt-GOA. We found some striking biases, stemming from the fact that a small fraction of articles that describe 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 mostly less

Only for certain ontologies though, right?

70 informative than those provided by low-throughput experiments; 2) annotations from
 71 high throughput experiments bias the annotations towards a limited number of functions,
 72 and, 3) many high-throughput experiments overlap in the proteins they annotate, and in
 73 the annotations assigned. Taken together, our findings offer a comprehensive picture of
 74 how the current protein function landscape is generated. Furthermore, due to the biases
 75 inherent in the current system of sequence annotations, this study serves as a caution to
 76 the producers and consumers of biological data from high-throughput experiments.

I really like the sentence, but are we overstating it?

Not so sure about the word "specific". Maybe something more like... Consequently, some studies now reveal certain types of functional characteristics for large groups of proteins as a result of the particular type of assay or assays used.

77 Methods and Results

78 Articles and Proteins

79 With the advent of high-throughput experiments it has become possible to conduct large-
 80 scale studies of protein functions. Consequently, some studies reveal very specific func-
 81 tional aspects of a large amount of proteins as a result of the particular type of assay or
 82 assays used. To understand the impact of large-scale studies on the corpus of experimen-
 83 tally annotated proteins, we looked at ~~the UniprotKB Gene Ontology (GO) annotation~~
 84 ~~files, or UniProt-GOA.~~ UniProt-GOA proteins are individually annotated by one or more
 85 GO terms using a procedure described in [6]. Briefly, this procedure consists of six steps
 86 which include sequence curation, sequence motif analyses, literature-based curation, recip-
 87 rocal BLAST [8] searches, attribution of all resources leading to the included findings, and
 88 quality assurance. If the annotation source is a research article, the attribution includes
 89 its PubMed ID. For each GO term associated with a protein, there is also an *evidence code*
 90 (EC) which is used to explain how the association between the protein and the GO term
 91 was made. Experimental evidence codes include such terms as: Inferred by Direct Assay
 92 (IDA) which indicates that "a direct assay was carried out to determine the function,

At entries in the UniprotKB Gene Ontology Annotation database, or UniProt-GOA.

Proteins in UniProt-GOA have been annotated with

93 process, or component indicated by the GO term” or Inferred from Physical Interaction
 94 (IPI) which “Covers physical interactions between the gene product of interest and an-
 95 other molecule.” (All EC definitions were taken from the GO site, geneontology.org).
 96 Computational evidence codes include terms such as *Inferred from Sequence or Structural*
 97 *Similarity* (ISS) and *Inferred from Sequence Orthology* (ISO). Although the evidence is
 98 non-experimental, the proteins annotated with these evidence codes are still assigned
 99 by a curator, rendering a level of human oversight. There are also ~~non~~-computational
 100 ~~and~~ non-experimental evidence codes, the most prevalent being *Inferred from Electronic*
 101 *Annotation* (IEA) which is “used for annotations that depend directly on computation
 102 or automated transfer of annotations from a database”. IEA evidence means that the
 103 annotation was not made or checked by a person. Different degrees of reliability are as-
 104 sociated with the ^{different} evidence codes, with experimental codes generally considered to be of
 105 higher reliability than non-experimental codes. However, the increase in the number of
 106 high-throughput experiments used to determine protein functions may introduce biases
 107 into experimental protein annotations, due to the inherent ^{experimental} capabilities and limitations
 108 of high-throughput assays. To test the hypothesis that such biases exist, and to study
 109 their extent if they do, we compiled the details of all experimentally-annotated proteins
 110 in UniProtKB. This included all proteins whose GO annotations have the GO experimen-
 111 tal evidence codes EXP, IDA, IPI, IMP, IGI, IEP. We first examined the distribution of
 112 articles by the number of proteins they annotate. As can be seen in **Figure 1**, the distri-
 113 bution of the number of proteins annotated per article follows a power-law distribution.
 114 $f(x) = ax^k$. Using linear regression over the log values of the axes we obtained a fit with
 115 $p < 1.18 \times 10^8$ and $R^2 = -0.72$. We therefore conclude that there is indeed a substantial
 116 bias in experimental annotations, in which there are few articles that annotate a large
 117 number of proteins.

Switch terms,
 use "but"
 instead of
 "and"

Alpha order?

Should be
 neg, i hope'

Should probably mention log odds again here.


118 To better understand the consequences of such a distribution, we divided the anno-
 119 tating articles into four cohorts, based on the number of proteins each article annotates.
 120 *Single-throughput* articles are those articles that annotate only one protein; *low through-*
 121 *put* articles annotate 2-9 proteins; *moderate throughput* articles annotate 10-99 proteins
 122 and *high throughput* articles annotate over 99 proteins. The results are shown in Table 1.
 123 The most striking finding is that high throughput articles are responsible for 25% of the
 124 annotations in Uniprot-GOA, even though they comprise 0.08% of the articles. 96% of
 125 the articles are single-throughput ~~and~~ ^{Or} low throughput, however those annotate only 53%
 126 of the proteins in Uniprot-GOA. So while moderate throughput and high-throughput ex-
 127 periments account for almost half of the annotations in Uniprot-GOA, they comprise only
 128 4% of the experiments published.

129 What typifies high-throughput articles? Also, how may the log-odds distribution
 130 bias what we understand of the protein function universe? To answer these questions,
 131 we examined different aspects of the annotations in the four article cohorts. Also, we
 132 examined in higher detail the top 50 annotating articles. (Overall, 62 articles in our

? Only 62 papers in exp uniprot-go
 are over 100? Just trying to make
 sure i understood the sentence...
 133 study annotated more than 100 proteins). An initial characterization of the top 50 high-
 134 throughput articles is shown in Table???. As can be seen, all of the articles are specific
 135 to a single species (typically a model organism) and assay that is used to annotate the
 136 proteins in that organism. Since a single assay was used, then typically only one ontology
 137 (MFO, BPO or CCO) was used for annotation. For some species this means that a
 138 single functional aspect (MFO, BPO or CCO) of a species will be dominated by a single
 139 experiment. ^{Category (?)}

?? Not sure what this means

140 Term frequency bias

141 To see how much a single species– and method– specific high-throughput assay affects
 142 the entire annotation of a species, we examined the relative contribution of the top-50
 143 articles to the entire corpus of experimentally annotated protein in each species. All the
 144 species found in the top-50 articles were either common model organisms or human. For
 145 each species, we looked at the five most frequent terms in the top 50 annotating articles.
 146 We then examined the contribution of this term by the top 50 articles to the general
 147 annotations of that species. The *contribution* is the number of annotations  by any given With (?)
 148 GO term in the top 50 articles divided by the number of annotations by that GO term in
 149 all of Uniprot-go ~~UniProtKB~~. For example, as seen in Figure 4 in *D. melanogaster* 88% of the usage
 150 of “precatalytic spliceosome” is contributed by the top-50 articles.

151 For most organisms in the top-50 articles, the annotations were within the cellular
 152 component ontology. The exceptions are *D. melanogaster* and *C. elegans* where the
 153 dominant terms were from the Biological Process ontology, and in mouse, where “protein
 154 binding” and “identical protein binding” are from the Molecular Function Ontology. *D.*
 155 *melanogaster*’s annotation for the top terms is dominated (over 50% contribution) by the
 156 top-50 articles.

157 The term frequency bias described here can be viewed more broadly within the ontol-
 158 ogy bias. The proteins annotated by the cohorts of single-protein articles, low-throughput
 159 articles, and moderate throughput articles have similar ratios of the fraction of proteins
 160 annotated. Twenty-two to twenty-six percent of assigned terms are in the Molecular
 161 Function Ontology, and 51-57% are in the Biological Process Ontology and the remaining
 162 17-25% are in the Cellular Component ontology. These ratios change dramatically with
 163 high-throughput articles (over 99 terms per article). In the high-throughput articles, only
 164 5% of assigned terms are in the Molecular Function Ontology, 38% in the Biological Pro-

165 cess Ontology and 57% in the Cellular Compartment Ontology, ostensibly due to a lack of
 166 high-throughput assays that can be used for generating annotations using the Molecular
 167 Function Ontology.

168 Reannotation Bias

Are we ok with bias here? Too negative sounding?

169 Another type of annotation bias is that of protein re-annotation. How many of the top-50
 170 articles actually re-annotate the same set of proteins? And how much of an agreement is
 171 there between different experiments? To investigate the extent of repetitive annotations in
 172 different articles, we clustered all the proteins annotated by the top-50 articles using CD-
 173 HIT [?] at 100% sequence identity. We then examined the number of clusters containing
 174 100% identical sequences per model species. The product of the number of proteins
 175 divided by the number of clusters is the redundancy percentage. For example, if each of
 176 the top-50 articles annotating the proteins in a given species annotated the same protein
 177 set, the redundancy percentage would be 100%. The results of the reannotation bias
 178 analysis are shown in Figure 2 and in Table 3. As can be seen, the highest redundancy
 179 (65%) is in the 12 articles annotating *C. elegans*.

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

186 Given a protein P , let G be the GO-terms g_1, g_2, \dots, g_m that annotate that protein in
 187 all top-50 articles for an ontology $O \in \{BPO, MFO, CCO\}$. The count of each of these
 188 go terms per protein per ontology is n_1, n_2, \dots, n_m with n_i being the number of times GO

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

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:

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

For example, the protein Oleate activated transcription factor 3 (UniProtID: P36023)

in *S. cerevisiae* is annotated by three articles in the CC ontology: **Probably write it out here not mixed abbreviate.**

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.

Says 2 articles in the figure legend.

Quantifying annotation information

A common assumption holds that while high-throughput experiments do annotate more protein functions than low-throughput experiments, the former also tend to be more

206 shallow in the predictions they provide. The information provided, for example, by a
 207 large-scale protein binding assay will only tell us if two proteins are binding, but will
 208 not reveal whether that binding is specific, will not provide an exact K_{bind} , will not say
 209 under what conditions binding takes place, or whether there is any enzymatic reaction
 210 or signal-transduction involved. Having on hand data from experiments with different
 211 “throughputness” levels, we set out to investigate whether there is a difference in the
 212 information provided by high-throughput experiments vs. low-throughput ones. To an-
 213 swer this question, we first had to quantify the information given by GO terms. One
 214 way to do so, is to use the depth of the term in the ontology: the term “enzyme activity”
 215 would be less informative than “dehalogenase” and the latter will be less informative than
 216 “haloalkane dehalogenase”. We therefore counted edges from the ontology root term to
 217 the GO-term to determine term information. The larger the number of edges, the more
 218 specific – and therefore informative – the annotation. In cases where several paths lead
 219 from the root to the examined GO-term, we used the minimal path. We did so for all the
 220 annotating articles split into groups by the number of proteins each article annotates.

221 Edge counting provides a measure of term-specificity. It is, however, imperfect. The
 222 reason is that different areas of the GO DAG have different connectivities, and terms
 223 may have different depths unrelated to the intuitive specificity of a term. For example
 224 “high-affinity Tryptophan transporter”, (GO:0005300) is 14 terms deep, while “anticoag-
 225 ulant”, (GO:0008435) is only three terms deep. For this reason, information content, the
 226 logarithm of the inverse of the GO term frequency in the corpus is generally accepted as a
 227 measure of GO-term information content [?]. To account for the possible bias created by
 228 the GO-DAG structure, we also used the log-frequency of the terms in the experimentally
 229 annotated proteins in Uniprot-GOA. However, it should be noted that the log-frequency
 230 measure is also imperfect because, as we see throughout this study, a GO-term’s frequency

Just to be
 clear, I
 suggest
 including
 level
 numbers
 here

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 in this study.

The results of both analyses are shown in Figure ?? 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

Have (?)



There are two complementary ways by which we come to knowledge about a protein's function. The twenty 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 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

255 terms than “Inferred by Direct Assay”: ECO also conveys which assay was used, i.e.
 256 “microscopy”, “RNA interference” etc. In addition to evidence terms, the ECO ontology
 257 provides *assertion terms* in which the nature of the assay is given. For example, an
 258 enzyme-linked immunosorbent assay (ELISA) provides quantitative protein data *in vitro*
 259 while an immunogold assay may provide the same information, and cellular localization
 260 information *in vivo*. It is therefore important to know both the assertion and the evi-
 261 dence to understand what sort of information may be gleaned from the assay. ~~However,~~
 262 to understand which types of assertions are made in the top-50 high throughput articles,
 263 we manually assigned Evidence Codes Ontology (ECO) assertion and evidence terms to
 264 the top-50 articles. The ECO ontology is more elaborate than the evidence codes used
 265 by Uniprot-GOA. Although there are plans to insert ECO terms into Uniprot-GOA in
 266 the near future, those will probably not be done manually for proteins already existing in
 267 Uniprot-GOA, but by automatic mapping EC terms to ECO ontology terms using a pre-
 268 set table (Rachael Huntley, Chris Mungall and Tony Sawford, personal communication).
 269 Thus, the ECO-based annotations we provide here to the top 50 articles is probably more
 270 informative than a future annotation may provide. Necessary?

271 The results are shown in Figure ?? and in Table 5.

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

278 The most used experimental assertion term was ECO:000160 “protein separation fol-
 279 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 limited. Second, there is the small number of terms used, when considering the large number of proteins that are being annotated. Third is the general term 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, *C. elegans* and *A. thaliana* studies comprise 36 (72%) of the top-50 annotating articles.

The most frequent experiment performed is cell fractionation and mass-spectrometry to assign a Cellular Component ontology terms (citations). 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, mass-spectrometry is the most common method used to localize proteins in subcellular compartments. A notable exception to the use of MS for protein localization is in the top annotating article [9] 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” [9]. 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 used** 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

Cut? How is it different from the previous sentence?

Study (?)

Experiment was

328 post-embryonic development [?]. Some studies focused on mitotic spindle assembly [10],
 329 lipid storage [10] and endocytic traffic [10]. One study used RNAi and MP to identify
 330 mitochondrial protein localization [11].

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

342 **Information Capture and Scope of GO**

343 We have discussed the information loss that is characteristic of high-throughput experi-
 344 ments, due to the nature of these experiments. However, another reason for information
 345 loss is the inability to capture certain types of information using the Gene Ontology. GO
 346 is knowingly limited to three aspects (MF, BP and CC) of biological function, which are
 347 assigned per protein. However, other aspects of function may emerge from experiments
 348 that cannot be captured by GO. Of note is the study mentioned earlier, “Proteome survey
 349 reveals modularity of the yeast cell machinery” [9]. In this study, the information pro-
 350 duced was primarily of protein complexes, which proteins are binding which proteins, and
 351 the relationship to cellular compartmentalization and biological networks. At the same

time, the only GO-term captured in the curation of this study was “protein binding”. Some, but not all of this information can be captured more specifically using the children of the term “protein binding”, but such a process is arguably laborious by manual curation of a high throughput article. Furthermore, the main information conveyed by this article, namely the types of protein complexes discovered and how they relate to cellular networks, is outside the scope of GO. It is important to realize that while high-throughput experiments do convey less information per protein within the functional scope as defined by GO, they still convey composite information such as possible pathway mappings – information which needs to be captured into annotation databases by means other than GO. In the example above, the information can be captured by a protein interaction database, but not by GO annotation.

I think we are possibly missing a paragraph on identical vs non identical annotations for 100% id sequences. What do you think?

So binding partners aren't supposed to be captured in go?

Conclusions

Taken together, the annotation biases noted in this study affect our understanding of protein function space. This, in turn, affects our ability to properly understand the connection between predictors of protein function and the actual function – the hallmark of computational function annotation. As a dramatic example, during the Critical Assessment of Function Annotation experiment (Radivojac *et al* in review) we have noticed that about 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 unknown proteins. If the training set for these programs has included a disproportionate number of annotations by high-throughput experiments, the results these programs provide will be strongly biased towards a few frequent and shallow GO-terms.

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 protein to determine its function by sequence similarity. If a target protein is tagged as annotated by a high throughput assay, it would be removed from the search if asked to do so by the user. This filtering can also be done by assay type, number of proteins annotated per experiment, or a combination of the above. This requires that GO-annotated proteins

should also be annotated with assertion codes in addition to the evidence codes and GO term-codes; but given the large volume of data in UniprotKB it is hard to expect such massive reannotation with assertion terms undertaken. (Any other ideas?)

We call upon the communities of annotators, computational biologists and experimental biologists to be mindful of the phenomenon of the biases described in this study, and to work together to understand its implications and mitigate its impact.

Isn't that happening, though, with goa moving to eco?

Things that could be recommended...

1. Adjust/ be careful of composition of training sets.
2. Use caution when transferring annotations from highthroughput sequences (because of unspecificity?)
3. Refrain from annotating sequence dbs with highthroughput paper info.
4. Use assertions/eco
5. Support standards & practices in publications to make finding function in a paper "easy" so that the temptation won't always be to add highthroughput papers because you "get more bang for your buck"

399 Note on Methods

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

404 Acknowledgments

405 We thank Predrag Radivojac, members of the Friedberg and Babbitt labs for insightful
406 discussions. This research was funded, in part by NSF / ABI XXXXXXXX-XXX award to
407 IF and XXXXXXXXXXXXX to PCB.

408 References

409 References

- 410 1. Friedberg I (2006) Automated protein function prediction—the genomic challenge.
411 Brief Bioinform 7: 225–242.
- 412 2. Erdin S, Lisewski AM, Lichtarge O (2011) Protein function prediction: towards
413 integration of similarity metrics. Current Opinion in Structural Biology 21: 180 -
414 188.
- 415 3. Rentzsch R, Orengo CA (2009) Protein function prediction the power of multiplic-
416 ity. Trends in Biotechnology 27: 210 - 219.

- 417 4. Sboner A, Mu X, Greenbaum D, Auerbach R, Gerstein M (2011) The real cost of
418 sequencing: higher than you think! *Genome Biology* 12: 125+.
- 419 5. Schnoes AM, Brown SD, Dodevski I, Babbitt PC (2009) Annotation error in public
420 databases: Misannotation of molecular function in enzyme superfamilies. *PLoS*
421 *Comput Biol* 5: e1000605+.
- 422 6. Dimmer EC, Huntley RP, Alam-Faruque Y, Sawford T, O'Donovan C, et al. (2012)
423 The uniprot-go annotation database in 2011. *Nucleic Acids Research* 40: D565–
424 D570.
- 425 7. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, et al. (2000) Gene ontology:
426 tool for the unification of biology. *Nature Genetics* 25: 25–29.
- 427 8. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, et al. (1997) Gapped
428 blast and psi-blast: a new generation of protein database search programs. *Nucleic*
429 *acids research* 25: 3389–3402.
- 430 9. Barbe L, Lundberg E, Oksvold P, Stenius A, Lewin E, et al. (2008) Toward a
431 confocal subcellular atlas of the human proteome. *Mol Cell Proteomics* 7: 499–
432 508.
- 433 10. Goshima G, Wollman R, Goodwin SS, Zhang N, Scholey JM, et al. (2007) Genes
434 required for mitotic spindle assembly in *Drosophila* S2 cells. *Science* 316: 417–421.
- 435 11. Hughes JR, Meireles AM, Fisher KH, Garcia A, Antrobus PR, et al. (2008) A
436 microtubule interactome: complexes with roles in cell cycle and mitosis. *PLoS Biol*
437 6: e98.

- 438 12. Matsuyama A, Arai R, Yashiroda Y, Shirai A, Kamata A, et al. (2006) ORFeome
439 cloning and global analysis of protein localization in the fission yeast *Schizosaccha-*
440 *romyces pombe*. *Nat Biotechnol* 24: 841–847.
- 441 13. Pagliarini DJ, Calvo SE, Chang B, Sheth SA, Vafai SB, et al. (2008) A mitochon-
442 drial protein compendium elucidates complex I disease biology. *Cell* 134: 112–123.
- 443 14. Simmer F, Moorman C, van der Linden AM, Kuijk E, van den Berghe PV, et al.
444 (2003) Genome-wide RNAi of *C. elegans* using the hypersensitive *rrf-3* strain reveals
445 novel gene functions. *PLoS Biol* 1: E12.
- 446 15. Huh WK, Falvo JV, Gerke LC, Carroll AS, Howson RW, et al. (2003) Global
447 analysis of protein localization in budding yeast. *Nature* 425: 686–691.
- 448 16. Zybaylov B, Rutschow H, Friso G, Rudella A, Emanuelsson O, et al. (2008) Sorting
449 signals, N-terminal modifications and abundance of the chloroplast proteome. *PLoS*
450 *ONE* 3: e1994.
- 451 17. Sonnichsen B, Koski LB, Walsh A, Marschall P, Neumann B, et al. (2005) Full-
452 genome RNAi profiling of early embryogenesis in *Caenorhabditis elegans*. *Nature*
453 434: 462–469.
- 454 18. Mootha VK, Bunkenborg J, Olsen JV, Hjerrild M, Wisniewski JR, et al. (2003)
455 Integrated analysis of protein composition, tissue diversity, and gene regulation in
456 mouse mitochondria. *Cell* 115: 629–640.
- 457 19. Benschop JJ, Mohammed S, O’Flaherty M, Heck AJ, Slijper M, et al. (2007) Quan-
458 titative phosphoproteomics of early elicitor signaling in *Arabidopsis*. *Mol Cell Pro-*
459 *teomics* 6: 1198–1214.

- 460 20. Kamath RS, Fraser AG, Dong Y, Poulin G, Durbin R, et al. (2003) Systematic
461 functional analysis of the *Caenorhabditis elegans* genome using RNAi. *Nature* 421:
462 231–237.
- 463 21. Mawuenyega KG, Forst CV, Dobos KM, Belisle JT, Chen J, et al. (2005) *My-*
464 *cobacterium tuberculosis* functional network analysis by global subcellular protein
465 profiling. *Mol Biol Cell* 16: 396–404.
- 466 22. Ito J, Batth TS, Petzold CJ, Redding-Johanson AM, Mukhopadhyay A, et al. (2011)
467 Analysis of the *Arabidopsis* cytosolic proteome highlights subcellular partitioning
468 of central plant metabolism. *J Proteome Res* 10: 1571–1582.
- 469 23. Rual JF, Ceron J, Koreth J, Hao T, Nicot AS, et al. (2004) Toward improving
470 *Caenorhabditis elegans* phenome mapping with an ORFeome-based RNAi library.
471 *Genome Res* 14: 2162–2168.
- 472 24. Reinders J, Zahedi RP, Pfanner N, Meisinger C, Sickmann A (2006) Toward the
473 complete yeast mitochondrial proteome: multidimensional separation techniques
474 for mitochondrial proteomics. *J Proteome Res* 5: 1543–1554.
- 475 25. Fernandez-Calvino L, Faulkner C, Walshaw J, Saalbach G, Bayer E, et al. (2011)
476 *Arabidopsis* plasmodesmal proteome. *PLoS ONE* 6: e18880.
- 477 26. Gu S, Chen J, Dobos KM, Bradbury EM, Belisle JT, et al. (2003) Comprehensive
478 proteomic profiling of the membrane constituents of a *Mycobacterium tuberculosis*
479 strain. *Mol Cell Proteomics* 2: 1284–1296.
- 480 27. Ferro M, Brugiere S, Salvi D, Seigneurin-Berny D, Court M, et al. (2010)
481 ATCHLORO, a comprehensive chloroplast proteome database with subplastidial

- 482 localization and curated information on envelope proteins. *Mol Cell Proteomics* 9:
483 1063–1084.
- 484 28. Kleffmann T, Russenberger D, von Zychlinski A, Christopher W, Sjolander K, et al.
485 (2004) The *Arabidopsis thaliana* chloroplast proteome reveals pathway abundance
486 and novel protein functions. *Curr Biol* 14: 354–362.
- 487 29. Sassetti CM, Boyd DH, Rubin EJ (2003) Genes required for mycobacterial growth
488 defined by high density mutagenesis. *Mol Microbiol* 48: 77–84.
- 489 30. Balklava Z, Pant S, Fares H, Grant BD (2007) Genome-wide analysis identifies
490 a general requirement for polarity proteins in endocytic traffic. *Nat Cell Biol* 9:
491 1066–1073.
- 492 31. Mitra SK, Gantt JA, Ruby JF, Clouse SD, Goshe MB (2007) Membrane proteomic
493 analysis of *Arabidopsis thaliana* using alternative solubilization techniques. *J Pro-*
494 *teome Res* 6: 1933–1950.
- 495 32. Maeda I, Kohara Y, Yamamoto M, Sugimoto A (2001) Large-scale analysis of gene
496 function in *Caenorhabditis elegans* by high-throughput RNAi. *Curr Biol* 11: 171–
497 176.
- 498 33. Ceron J, Rual JF, Chandra A, Dupuy D, Vidal M, et al. (2007) Large-scale RNAi
499 screens identify novel genes that interact with the *C. elegans* retinoblastoma path-
500 way as well as splicing-related components with synMuv B activity. *BMC Dev Biol*
501 7: 30.
- 502 34. Sickmann A, Reinders J, Wagner Y, Joppich C, Zahedi R, et al. (2003) The pro-
503 teome of *Saccharomyces cerevisiae* mitochondria. *Proc Natl Acad Sci USA* 100:
504 13207–13212.

- 505 35. Gavin AC, Aloy P, Grandi P, Krause R, Boesche M, et al. (2006) Proteome survey
506 reveals modularity of the yeast cell machinery. *Nature* 440: 631–636.
- 507 36. Green RA, Kao HL, Audhya A, Arur S, Mayers JR, et al. (2011) A high-resolution
508 *C. elegans* essential gene network based on phenotypic profiling of a complex tissue.
509 *Cell* 145: 470–482.
- 510 37. Simpson JC, Wellenreuther R, Poustka A, Pepperkok R, Wiemann S (2000) Sys-
511 tematic subcellular localization of novel proteins identified by large-scale cDNA
512 sequencing. *EMBO Rep* 1: 287–292.
- 513 38. Marmagne A, Ferro M, Meinnel T, Bruley C, Kuhn L, et al. (2007) A high content
514 in lipid-modified peripheral proteins and integral receptor kinases features in the
515 arabidopsis plasma membrane proteome. *Mol Cell Proteomics* 6: 1980–1996.
- 516 39. Dunkley TP, Hester S, Shadforth IP, Runions J, Weimar T, et al. (2006) Mapping
517 the Arabidopsis organelle proteome. *Proc Natl Acad Sci USA* 103: 6518–6523.
- 518 40. Jaquinod M, Villiers F, Kieffer-Jaquinod S, Hugouvieux V, Bruley C, et al. (2007)
519 A proteomics dissection of Arabidopsis thaliana vacuoles isolated from cell culture.
520 *Mol Cell Proteomics* 6: 394–412.
- 521 41. Heazlewood JL, Tonti-Filippini JS, Gout AM, Day DA, Whelan J, et al. (2004) Ex-
522 perimental analysis of the Arabidopsis mitochondrial proteome highlights signaling
523 and regulatory components, provides assessment of targeting prediction programs,
524 and indicates plant-specific mitochondrial proteins. *Plant Cell* 16: 241–256.
- 525 42. Ashrafi K, Chang FY, Watts JL, Fraser AG, Kamath RS, et al. (2003) Genome-
526 wide RNAi analysis of *Caenorhabditis elegans* fat regulatory genes. *Nature* 421:
527 268–272.

- 528 43. Piano F, Schetter AJ, Morton DG, Gunsalus KC, Reinke V, et al. (2002) Gene
529 clustering based on RNAi phenotypes of ovary-enriched genes in *C. elegans*. *Curr*
530 *Biol* 12: 1959–1964.
- 531 44. Carter C, Pan S, Zouhar J, Avila EL, Girke T, et al. (2004) The vegetative vacuole
532 proteome of *Arabidopsis thaliana* reveals predicted and unexpected proteins. *Plant*
533 *Cell* 16: 3285–3303.
- 534 45. Da Cruz S, Xenarios I, Langridge J, Vilbois F, Parone PA, et al. (2003) Proteomic
535 analysis of the mouse liver mitochondrial inner membrane. *J Biol Chem* 278: 41566–
536 41571.
- 537 46. Rual JF, Venkatesan K, Hao T, Hirozane-Kishikawa T, Dricot A, et al. (2005)
538 Towards a proteome-scale map of the human protein-protein interaction network.
539 *Nature* 437: 1173–1178.
- 540 47. Bakthavatsalam D, Gomer RH (2010) The secreted proteome profile of developing
541 *Dictyostelium discoideum* cells. *Proteomics* 10: 2556–2559.
- 542 48. Froehlich JE, Wilkerson CG, Ray WK, McAndrew RS, Osteryoung KW, et al.
543 (2003) Proteomic study of the *Arabidopsis thaliana* chloroplastic envelope mem-
544 brane utilizing alternatives to traditional two-dimensional electrophoresis. *J Pro-*
545 *teome Res* 2: 413–425.
- 546 49. Stroschein-Stevenson SL, Foley E, O’Farrell PH, Johnson AD (2006) Identification
547 of *Drosophila* gene products required for phagocytosis of *Candida albicans*. *PLoS*
548 *Biol* 4: e4.

- 549 50. Rutschow H, Ytterberg AJ, Friso G, Nilsson R, van Wijk KJ (2008) Quantitative
550 proteomics of a chloroplast SRP54 sorting mutant and its genetic interactions with
551 CLPC1 in Arabidopsis. *Plant Physiol* 148: 156–175.
- 552 51. Kumar A, Agarwal S, Heyman JA, Matson S, Heidtman M, et al. (2002) Subcellular
553 localization of the yeast proteome. *Genes Dev* 16: 707–719.
- 554 52. Fraser AG, Kamath RS, Zipperlen P, Martinez-Campos M, Sohrmann M, et al.
555 (2000) Functional genomic analysis of *C. elegans* chromosome I by systematic RNA
556 interference. *Nature* 408: 325–330.
- 557 53. Gonczy P, Echeverri C, Oegema K, Coulson A, Jones SJ, et al. (2000) Functional
558 genomic analysis of cell division in *C. elegans* using RNAi of genes on chromosome
559 III. *Nature* 408: 331–336.
- 560 54. Suzuki H, Fukunishi Y, Kagawa I, Saito R, Oda H, et al. (2001) Protein-protein
561 interaction panel using mouse full-length cDNAs. *Genome Res* 11: 1758–1765.
- 562 55. Sarry JE, Kuhn L, Ducruix C, Lafaye A, Junot C, et al. (2006) The early re-
563 sponses of *Arabidopsis thaliana* cells to cadmium exposure explored by protein and
564 metabolite profiling analyses. *Proteomics* 6: 2180–2198.
- 565 56. Chen D, Toone WM, Mata J, Lyne R, Burns G, et al. (2003) Global transcriptional
566 responses of fission yeast to environmental stress. *Mol Biol Cell* 14: 214–229.
- 567 57. Herold N, Will CL, Wolf E, Kastner B, Urlaub H, et al. (2009) Conservation of the
568 protein composition and electron microscopy structure of *Drosophila melanogaster*
569 and human spliceosomal complexes. *Mol Cell Biol* 29: 281–301.

- 570 58. Bayer EM, Bottrill AR, Walshaw J, Vigouroux M, Naldrett MJ, et al. (2006) Ara-
571 bidopsis cell wall proteome defined using multidimensional protein identification
572 technology. *Proteomics* 6: 301–311.

573 Figures

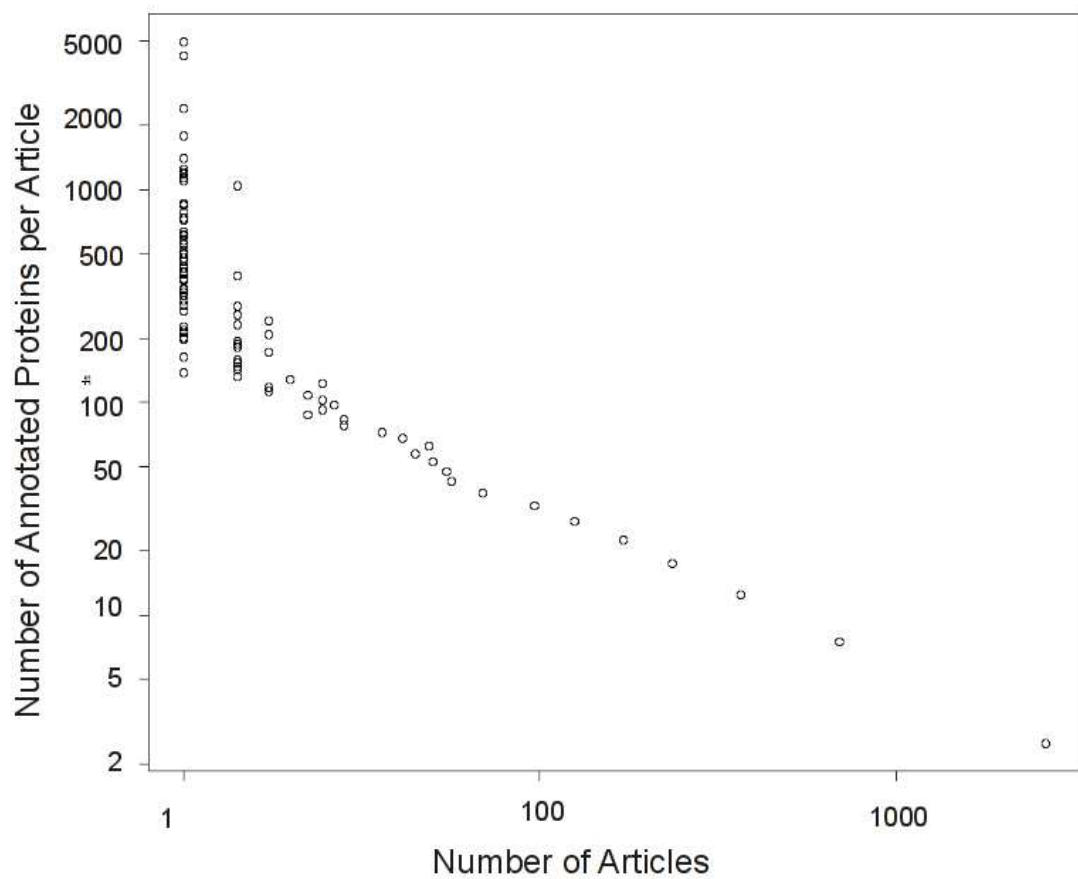


Figure 1. Distribution of the number of proteins annotated per article.

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.

Not the same p value as in the paper

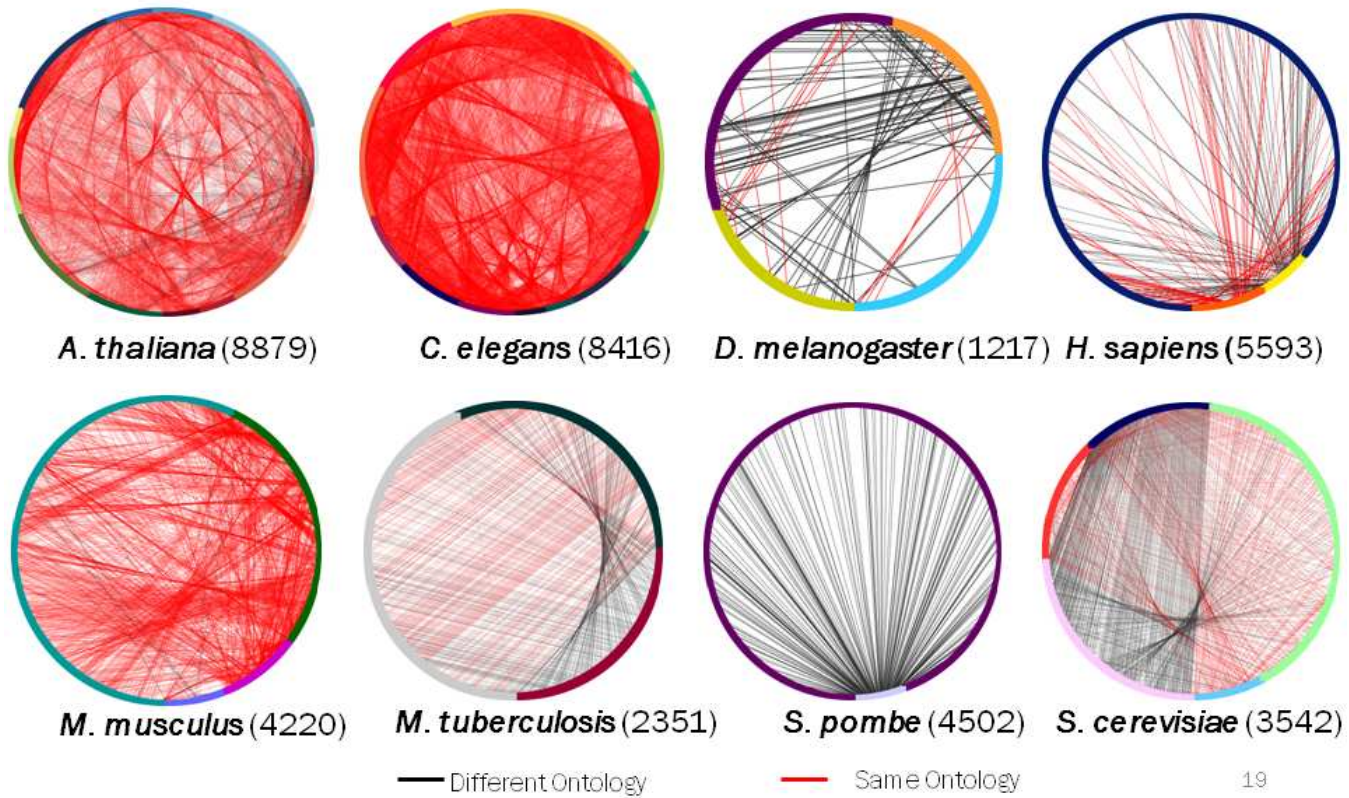


Figure 2. Redundancy in proteins described by the top-50 articles. A circle represents the sum total of articles annotating each organism. Each colored arch is composed of all the proteins in a single article. A line is drawn between any two points on the circle if the proteins they represent have 100% 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.

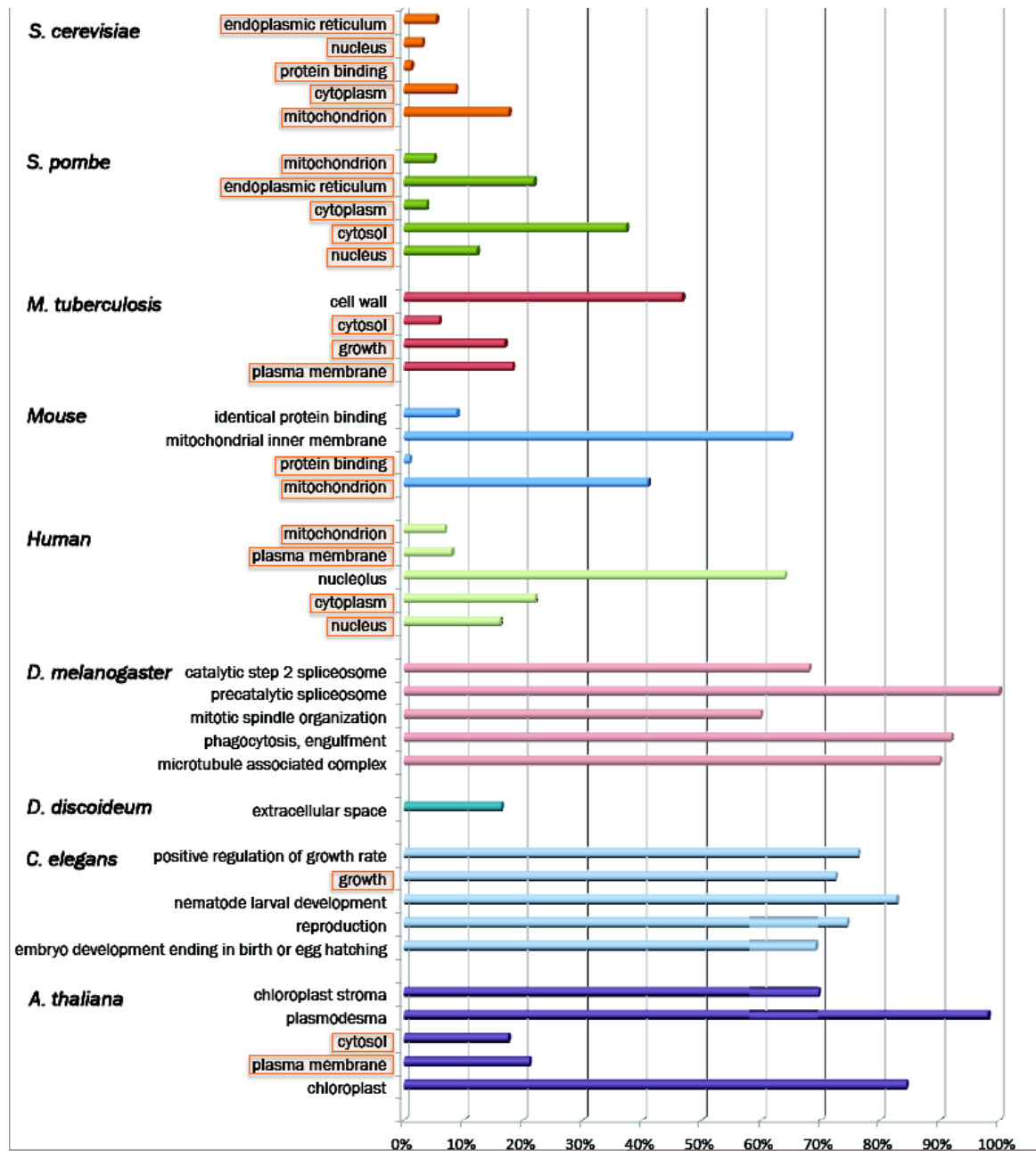


Figure 3. Relative contribution of top-50 articles to the annotation of major model organisms. The length of each bar represents the percentage of proteins annotated by the top-50 articles in a given organism by a given GO term.

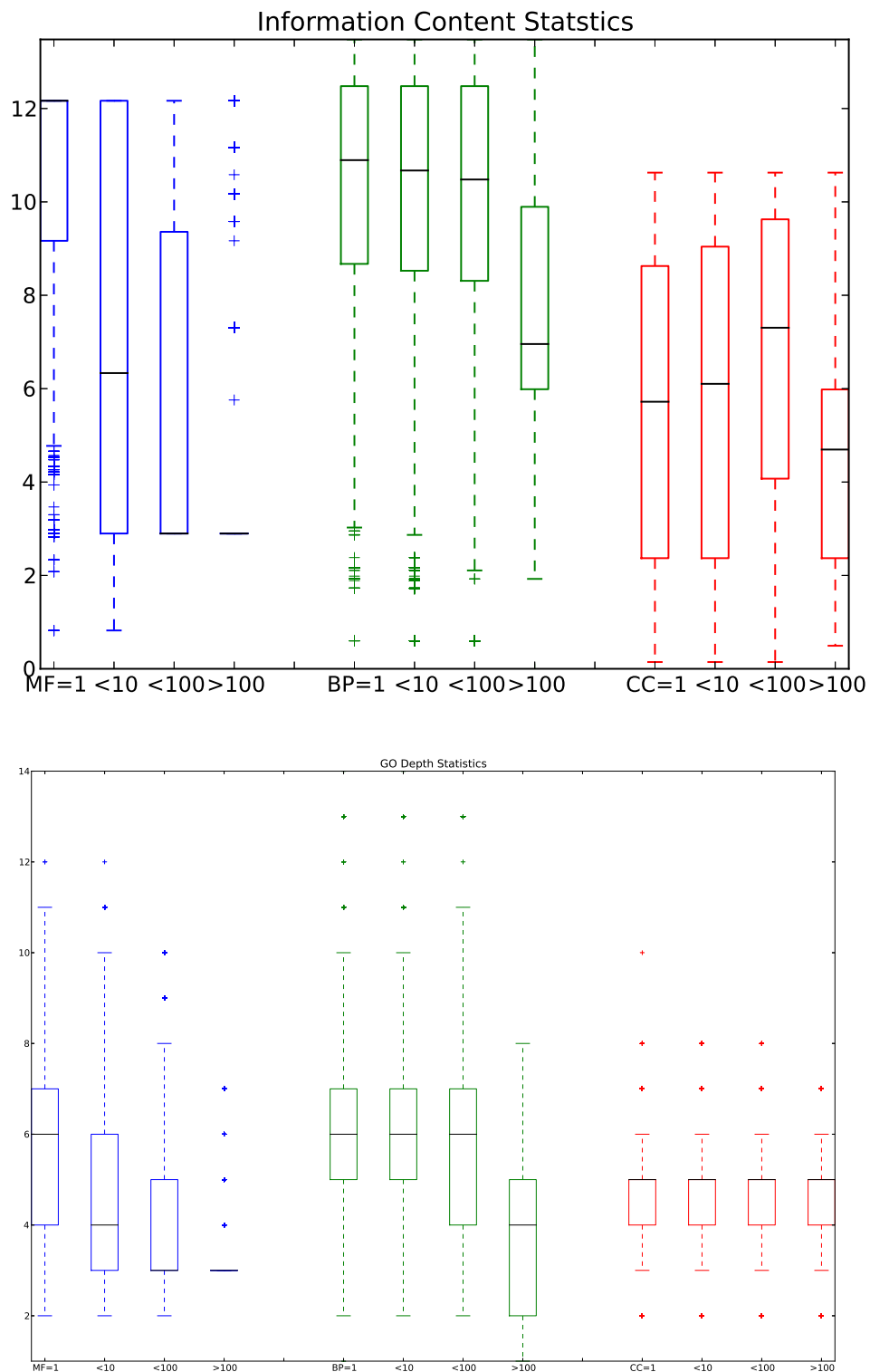


Figure 4. Information provided by articles depending on the number of proteins the articles annotate. Articles are grouped into cohorts: 1: one protein annotated by article; <10: more than 1, less than 10 annotated; <100: more than 10, less than 100 annotated; ≥ 100: more than 100 proteins annotated 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.

574 **Tables**

Table 1. 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	2668	62	76087
Percent of proteins annotated	16.56	37.11	21.19	25.13	100
Percent of annotating articles	54.09	42.32	3.51	0.08	100

Table caption

Table 2. Top 50 Annotating Articles

N	Proteins	Annotations	Species	ref.	MFO	BPO	CCO
1	4937	11050	<i>H. sapiens</i>	[9]	0	0	11050
2	4247	7046	<i>S. pombe</i>	[12]	0	0	7046
3	2412	2412	<i>H. sapiens</i>	[13]	0	0	2412
4	1791	5918	<i>C. elegans</i>	[14]	0	5918	0
5	1406	1863	<i>S. cerevisiae</i>	[15]	0	0	1863
6	1251	1251	<i>A. thaliana</i>	[16]	0	0	1251
7	1205	1476	<i>C. elegans</i>	[17]	0	1476	0
8	1186	1213	<i>M. musculus</i>	[18]	0	0	1213
9	1136	1136	<i>A. thaliana</i>	[19]	0	0	1136
10	1101	2269	<i>C. elegans</i>	[20]	0	2269	0
11	1043	1365	<i>M. tuberculosis</i>	[21]	0	0	1365
12	1041	1041	<i>A. thaliana</i>	[22]	0	0	1041
13	865	1533	<i>C. elegans</i>	[23]	0	1533	0
14	845	845	<i>S. cerevisiae</i>	[24]	0	0	845
15	784	784	<i>A. thaliana</i>	[25]	0	0	784
16	735	735	<i>M. tuberculosis</i>	[26]	0	0	735
17	724	882	<i>A. thaliana</i>	[27]	0	0	882
18	634	634	<i>A. thaliana</i>	[28]	0	0	634
19	613	613	Mycobacter sp.	[29]	0	613	0
20	607	661	<i>C. elegans</i>	[30]	0	659	2
21	577	577	<i>A. thaliana</i>	[31]	0	0	577

Continued on next page

N	Proteins	Annotations	Species	ref.	MFO	BPO	CCO
22	553	884	<i>C. elegans</i>	[32]	0	884	0
23	516	5972	<i>C. elegans</i>	[33]	0	5972	0
24	503	503	<i>S. cerevisiae</i>	[34]	0	0	503
25	498	638	<i>S. cerevisiae</i>	[35]	638	0	0
26	479	848	<i>C. elegans</i>	[36]	0	848	0
27	465	468	<i>H. sapiens</i>	[37]	0	0	468
28	436	436	<i>A. thaliana</i>	[38]	0	0	436
29	430	513	<i>A. thaliana</i>	[39]	0	0	513
30	413	456	<i>D. melanogaster</i>	[11]	0	39	417
31	401	401	<i>A. thaliana</i>	[40]	0	0	401
32	392	392	<i>A. thaliana</i>	[41]	0	0	392
33	392	639	<i>C. elegans</i>	[42]	0	639	0
34	383	917	<i>C. elegans</i>	[43]	0	917	0
35	380	380	<i>A. thaliana</i>	[44]	0	0	380
36	375	375	<i>M. musculus</i>	[45]	0	0	375
37	343	509	<i>H. sapiens</i>	[46]	509	0	0
38	338	338	Ddiscoideum	[47]	0	0	338
39	328	328	<i>A. thaliana</i>	[48]	0	0	328
40	319	329	<i>C. albicans</i>	[49]	1	328	0
41	305	312	<i>A. thaliana</i>	[50]	0	0	312
42	290	331	<i>S. cerevisiae</i>	[51]	0	0	331
43	285	761	<i>C. elegans</i>	[52]	0	761	0

Continued on next page

N	Proteins	Annotations	Species	ref.	MFO	BPO	CCO
44	283	499	<i>C. elegans</i>	[53]	0	499	0
45	266	433	<i>M. musculus</i>	[54]	433	0	0
46	260	260	<i>A. thaliana</i>	[55]	0	260	0
47	258	259	<i>S. pombe</i>	[56]	0	259	0
48	244	397	<i>D. melanogaster</i>	[10]	0	367	30
49	242	397	<i>D. melanogaster</i>	[57]	0	0	397
50	241	263	<i>A. thaliana</i>	[58]	0	0	263

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

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	Ontology	num prot	mean $k_{P,O}$	stdv	stderr	num articles	num t
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 are also provided. **num articles:** number of annotating articles **num terms** number of annotating terms. Annotations by less than 2 articles or two terms (or both) for the same protein/ontology combination have been omitted.

Different in paper text

Table 5. 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

Table caption