

1   **TITLE**

2   A computational framework for analysis of prey-prey associations in interaction proteomics identifies  
3   novel human protein-protein interactions and networks

4   **AUTHORS AND AFFILIATIONS**

5   Sudipto Saha<sup>1</sup>, Jean-Eudes Dazard<sup>1</sup>, Hua Xu<sup>1</sup> and Rob M. Ewing<sup>1,2\*</sup>

6   <sup>1</sup>Center for Proteomics and Bioinformatics, <sup>2</sup>Department of Genetics and Genome Sciences,  
7   Case Western Reserve University School of Medicine, Cleveland, Ohio 44106, USA

8  
9   \*To whom correspondence should be addressed: Rob M. Ewing, Ph.D., Center for Proteomics and  
10   Bioinformatics, Case Western Reserve University, 10900 Euclid Avenue, Cleveland, OH 44106; Tel:  
11   (216) 368-4380; Fax: (216) 368-6846, Email: [rob.ewing@case.edu](mailto:rob.ewing@case.edu)

12   **Abstract**

13   Large-scale protein-protein interaction datasets have been generated for several species including yeast  
14   and human and have enabled the identification, quantification and prediction of cellular molecular  
15   networks. Affinity purification-mass spectrometry (AP-MS) is the preeminent methodology for large-  
16   scale analysis of protein complexes, performed by immunopurifying a specific 'bait' protein and its  
17   associated 'prey' proteins. The analysis and interpretation of AP-MS datasets is however, not  
18   straightforward. In addition, although yeast AP-MS datasets are relatively comprehensive, current  
19   human AP-MS datasets only sparsely cover the human interactome. Here we develop a framework for  
20   analysis of AP-MS datasets that addresses the issues of noise, missing data and sparsity of coverage in  
21   the context of a current, real world human AP-MS dataset. Our goal is to extend and increase the  
22   density of the known human interactome by integrating bait-prey and co-complexed preys (prey-prey

associations) into networks. Our framework incorporates a score for each identified protein, as well as elements of signal processing to improve the confidence of identified protein-protein interactions. We identify many protein networks enriched in known biological processes and functions. In addition, we show that integrated bait-prey and prey-prey interactions can be used to refine network topology and extend known protein networks.

## **Keywords**

protein-protein interaction network; affinity purification mass-spectrometry; interactome

## **Introduction**

Interactions between proteins and the protein complexes and networks that these interactions form, are fundamental units of biological organization that mediate most cellular processes. Understanding the topology of protein interaction networks is therefore a primary goal of systems biology. The principal techniques for large-scale analyses of protein interactions are yeast two-hybrid and affinity-purification mass-spectrometry (AP-MS). These two approaches provide complementary views of the protein interactome<sup>1</sup>; the yeast-two-hybrid assay identifies binary interactions between pairs of proteins, whilst affinity-purification mass-spectrometry (AP-MS) identifies protein complexes associated with a given “bait” protein. AP-MS experiments have been conducted by purifying protein complexes using native antibodies or on a larger-scale by epitope tagging bait proteins and recovering associated protein complexes with antibodies directed against the epitope tag. Despite the power of AP-MS to map protein complexes, however, several technical hurdles exist, including the presence of non-specific interacting proteins in these large-scale datasets as well as lack of reproducibility and sparse coverage of the underlying networks.

The yeast protein interactome is by far the most well covered eukaryotic interactome with multiple studies and techniques contributing large-scale datasets<sup>2-5</sup>. With its larger size and greater complexity,

the human protein interactome has been mapped more selectively, with studies focusing on specific complexes or classes of protein<sup>6-9</sup>. Diverse computational methods have been developed for analysis of large-scale interaction proteomics datasets. For AP-MS datasets, the focus of these analysis methods has been the assignment of confidence scores to observed interactions and to distinguish specific from non-specific interactions<sup>6-7,10-14</sup>. In yeast, the socio-affinity index was developed to score a large-scale yeast interactome dataset with high coverage, and reciprocal AP-MS experiments<sup>10</sup>. In human datasets, with less dense coverage of the underlying networks, computational methods and scores have focused on distinguishing specific from non-specific interactions. For example, the normalized spectral abundance factor (NSAF) was used as a measure of abundance for each protein and the ratio of the vectors of counts for each protein between “control” and “bait” experiments used to distinguish specific from non-specific interactors<sup>6</sup>. The D-score, a metric combining total spectral count, reproducibility and overall frequency of prey proteins in AP-MS experiments into a single confidence score for bait-prey associations was developed for the analysis of human AP-MS data<sup>7</sup>. The D-score was applied to a large-scale study of human deubiquitinating enzymes and shown to outperform the socio-affinity index, the NSAF method and Z-score on the data in hand<sup>7</sup>. An alternative approach, using mixture modeling and Bayesian statistics named SAINT (Significance Analysis of Interactome) was developed<sup>13</sup> and applied to the analysis of phosphatase interaction networks<sup>15</sup>. The Decontaminator uses a Bayesian approach to model false-positive protein-protein interactions (PPIs) by comparing the score of a putative prey protein in induced vs control experiments<sup>14</sup>. Two broad models for interpreting AP-MS data in a network context have been proposed. The “spoke” model holds that the bait protein interacts with each of the identified “prey” proteins (the bait representing the center of a wheel with spokes connecting to each of the prey proteins) whereas the “matrix” model assumes that each of the identified proteins (bait and prey) interacts with each of the others in a given AP-MS experiment<sup>16</sup>. Although the matrix model will capture a higher proportion of

the underlying protein associations, this comes with the price of increased false positives. Therefore, although both models have merits, it is likely that a combination of spoke and matrix models represent underlying biological reality in most cases. Several scoring metrics for AP-MS data explicitly implement these concepts. The socio-affinity index is a summary score including spoke and matrix terms, as well as accounting for reciprocal bait-prey AP-MS experiments and computes the log ratio of protein co-occurrence given their observed frequencies<sup>10</sup>. The hypergeometric distribution has been used to calculate expected frequencies of co-occurrence for proteins using the matrix model in the large-scale yeast AP-MS datasets and found to be an effective means of identifying protein-protein interactions<sup>17</sup>. Other schemes that make use of co-occurrence frequencies have been proposed for resolving protein complexes using large-scale yeast AP-MS datasets<sup>18-19</sup> but with added refinements such as consideration of the variation in bait affinity when computing the results<sup>18</sup>.

Probabilistic approaches using small to medium scale human AP-MS datasets have also been developed<sup>6,20</sup>. These studies analyze datasets in which the coverage of given protein complexes by bait proteins is relatively high. Notably, these studies<sup>6,20</sup> make use of the quantitative features of mass-spectrometry data, rather than using binary co-occurrence frequencies. Measures of abundance such as spectral counts and other mass-spectrometry based confidence measures are very informative features of AP-MS data, since protein-protein affinities are likely to vary considerably. An emerging consensus of these diverse computational methods is that both bait-prey and prey-prey associations in AP-MS datasets can be used to identify protein-protein interactions. Utilization of prey-prey associations is most obvious in AP-MS datasets with dense bait coverage, so that co-occurrence profiles are well defined. However, it is less clear whether this concept may be applied to less dense AP-MS datasets such as larger-scale human AP-MS datasets<sup>7,12</sup>. Our goal in this work is to test the utility of mining prey-prey associations from large-scale, but less dense human AP-MS data.

We previously generated a large-scale human interaction proteomics dataset focused on 338 human bait



94 proteins, many of which are linked to human diseases<sup>12</sup>. This dataset is the largest human AP-MS  
95 dataset to date and mass-spectrometry experiments were performed in a uniform manner, making it a  
96 useful resource for development of data analysis techniques. In contrast to other human AP-MS  
97 datasets that focus on individual complexes or processes, our dataset represents a broad survey of  
98 human proteins and their interactions. In addition, the dataset represents several complexes with  
99 relatively high coverage (in terms of number of baits), such as the proteasome, and Eukaryotic  
100 Initiation Factors, and many other human protein complexes that are sparsely covered. Thus our  
101 motivation is to develop a method that can identify probable protein complexes from a heterogeneous  
102 (in terms of coverage or sampling of protein complexes) dataset. In our initial analysis of the dataset we  
103 focused on bait-prey interactions (spoke model) by using partial least squares to predict the  
104 reproducibility of each prey protein observation based upon a training set of highly-reproduced AP-MS  
105 experiments<sup>12</sup>. In the current work, we mine the dataset more comprehensively by identifying prey-  
106 prey associations based upon co-occurrence profiles. We show that these prey-prey interactions are a  
107 valuable source of protein interactions and that integration of bait-prey and prey-prey associations  
108 yields improved network models of protein complexes. Importantly, we show that prey-prey  
109 associations may be identified from less well covered interactomes, such as the current human  
110 interactome. Since comprehensive efforts to map the human protein interactome are still in their early  
111 states, methods to mine protein-protein interactions from incomplete AP-MS datasets will be important  
112 tools for the foreseeable future.

113 Our framework exploits the quantitative features of AP-MS data to assign confidence scores to each  
114 bait-prey observation. To improve signal-to-noise ratio, the sparse matrix of scores are then  
115 transformed to eigenvalues using singular value decomposition. A similarity measure between the  
116 profiles of each pair of prey proteins is then to detect potential prey-prey associations, which are the  
117 starting point for the clustering and reconstruction of protein complexes, which we validate by

reference to known complexes and annotations. Finally, we integrate bait-prey and prey-prey associations and show that these integrated networks have improved biological coherence. Importantly, we are able to directly compare bait-prey and prey-prey interactions in terms of their biological coherence and show that prey-prey interactions are a rich source of protein-protein associations. We illustrate our approach using selected specific protein networks and show how novel interactions can be identified.

The principal contribution of our work is to show how large-scale AP-MS datasets may be data-mined for identification of protein-protein associations. Our analysis pipeline is generically applicable to large-scale AP-MS datasets and we anticipate that as increasing volumes of AP-MS data are available, it can be applied and used to discover novel protein interactions and elucidate the topology of protein interaction networks.

## **Methods**

### *Data set*

Data used to develop our approach is principally derived from a previously described human AP-MS dataset, in which 832 single-step anti-FLAG immunoprecipitation experiments representing 384 human bait proteins (~50% of baits were replicated) identified 5269 distinct prey proteins<sup>12</sup>. The human bait proteins were selected based on association with diseases such as cancer and obesity. The data were generated as previously described<sup>12</sup>, except that all spectra were re-searched against an IPI human protein sequence database (version 3.31) (92012 sequences) using the Mascot mass-spectrometry search engine (version 2.1, Matrix Science; fixed modification: Carbamidomethyl (C), variable modification: Oxidation (M); peptide mass tolerance 2Da; fragment mass tolerance 0.4Da; missed cleavages:2). Peptide and protein identifications were processed through Peptide and Protein Prophet<sup>43</sup> and imported to LabKey server<sup>44</sup> (version 10.3) for data management. For subsequent steps of the

141 analysis, corresponding gene symbols, where available, were used to represent baits and preys.

#### 142 *Protein Identification False Discovery Rates*

143 To assess the protein identification false discovery rate, we used decoy database searches. Using the  
144 same Mascot search parameters as above, we searched all data against a concatenated decoy human IPI  
145 protein sequence database (version 3.31) and computed the mean false discovery rate across all  
146 searches (3.63%). We also searched the complete datasets using another search-engine, MassMatrix<sup>42</sup>  
147 (version 2.4.2, <http://www.massmatrix.net>) that provides the false discovery rate in terms of the %  
148 decoy hits when searched against concatenated decoy databases. The database and search parameters  
149 were the same as for the Mascot analysis, except for the following additional options: peptide length: 6-  
150 40 amino acid residues and score thresholds of 5.3 and 1.3 for the pp and pp<sub>tag</sub> scores respectively.  
151 Proteins identified using MassMatrix were cross-referenced to those identified using Mascot. Of the set  
152 of 34383 bait-prey associations from the original Mascot searches, and for which a D-score was  
153 computed (see below), 82% were also identified with MassMatrix. The mean false discovery rate  
154 across this set is 4.52%.

#### 155 *Bait-prey confidence score*

156 The D-score (Equation 1) is a previously described confidence score for AP-MS data that combines  
157 measures of abundance (spectral counts), specificity and reproducibility into a single score for each  
158 bait-prey observation<sup>7</sup>.

$$D_{ij} = \sqrt{\left( \frac{k}{\sum_{i=1}^k f_{ij}} \right)^p} X_{ij} \text{ for all } i, j \quad (\text{Eq. 1})$$

160

161 Where,  $k$ = total number of unique bait proteins;  $X_{ij}$ = total spectral counts for prey  $i$  from bait  $j$ ;  
162  $f_{ij}=\{1,0\}$ ;  $p$ = number of replicates runs in which the prey protein is present. This generated a complete  
163 bait-prey D-score matrix (**D**) of dimensions (5269 proteins x 384 baits).

#### 164 *Contaminant identification*

165 False-positive proteins may occur in AP-MS experiments for different reasons. First, there are many  
166 proteins that occur frequently in AP-MS experiments as a result of non-specific affinity. In the AP-MS  
167 dataset used in this study, a set of 200 control AP-MS experiments (using HEK293 cells with vector  
168 alone, i.e. no bait) provide a dataset for identifying 'control' proteins. Since the D-score accounts for  
169 protein frequency of occurrence, these proteins typically have low D-scores. D-scores for the top 100  
170 highly frequent control proteins in this study have median value equal to 3 (Supplementary Figure S1).  
171 These were removed along with any protein with D-score  $<3$  from the initial dataset (13688 bait-prey  
172 interactions with D-score  $\leq 3$  were removed). A second type of contaminant results from cross-  
173 contamination at the experimental level between samples. Since the experimental protocol used to  
174 generate data used here<sup>12</sup> used gel-based separation of proteins prior to mass-spectrometry, we  
175 identified proteins with D-score  $> 10$  co-occurring on the same gels or samples that were processed on  
176 the same date with different bait proteins as potential cross-contaminants and removed them from the  
177 dataset (1550 bait-prey interactions were removed using this criterion).

#### 178 *Latent Semantic Analysis (LSA)*

179 The D-score matrix (**D**) is necessarily sparse (1.25% non-zero values), since the frequency of most  
180 proteins is low (identified with few or unique baits). To identify relationships between proteins in this  
181 matrix, we made use of the principles of latent semantic analysis (LSA) which uses Singular Value  
182 Decomposition (SVD) to reduce matrix sparseness<sup>21</sup>. SVD was applied to the rectangular bait-prey

183 matrix (Equation 2), and a dimensional reduction of the singular values was performed after identifying  
184 the optimum  $k$ -value which determines the degree of reduction.

185

$$186 \quad \mathbf{D} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \approx \mathbf{D}_k = \mathbf{U}_k\mathbf{\Sigma}_k\mathbf{V}_k^T \quad (\text{Eq. 2})$$

187

188  $\mathbf{D}$  denotes the D-score bait-prey matrix ( $m \times n$ ).  $\mathbf{U}$  ( $m \times m$  matrix) and  $\mathbf{V}$  ( $n \times n$  matrix) are orthogonal  
189 matrices with unit-length columns (i.e.,  $\mathbf{U}\mathbf{U}^T = \mathbf{I}$  and  $\mathbf{V}^T\mathbf{V} = \mathbf{I}$ ) and  $\mathbf{\Sigma}$  is a diagonal matrix containing  
190 the ordered singular values of rank  $r = \min(m, n)$ .  $\mathbf{D}_k$  denotes the approximated latent semantic  
191 representation of the bait-prey matrix  $\mathbf{D}$  where  $k$  denotes the selected number of singular values used in  
192 the approximation. This generated a final bait-prey D-score matrix ( $\mathbf{D}_k$ ) of dimensions (2242 proteins x  
193 384 baits). LSA computations were performed using the Perl Data Language (PDL, <http://pdl.perl.org>).

194

#### 195 *Prey-prey similarity score*

196 For each bait protein, the bait vector is the vector of D-scores for all preys. Similarly, for each prey  
197 protein, the prey vector is the vector of all D-scores for all baits. LSA projects the vectors of original  
198 bait and prey vectors into lower dimensional semantic space. Transformed bait and prey vectors are the  
199 vectors of bait or prey eigenvalues. Cosine similarity (Equation 3) was used to compute pairwise  
200 similarities between prey-prey vectors from the bait-prey matrix. The similarity score for a pair of prey  
201 proteins a, and b is then the cosine similarity of their associated D-score vectors (0 indicates minimal  
202 similarity, 1 indicates maximal similarity). Hereafter the cosine similarity is referred to as the prey-prey  
203 score (PPS).

204

$$\text{sim}(a, b) \equiv \cos \alpha(a, b) \equiv \frac{a \cdot b}{\|a\| \|b\|} = \frac{\left( \sum_{i=1}^n a_i b_i \right)}{\left( \sqrt{\sum_{i=1}^n a_i^2} \sqrt{\sum_{i=1}^n b_i^2} \right)} \quad (\text{Eq. 3})$$

### *Analysis of known interactions*

Bait-prey and prey-prey interactions were compared to CORUM version 2.0<sup>22</sup> and BioGRID version 3.1.78<sup>23</sup> datasets. The complete set of CORUM mammalian protein complexes was used (2577 protein complexes; 4304 distinct gene Ids). Human protein interactions were extracted from BioGRID (402127 total interactions; 37910 human interactions). Data originating from our previous publication<sup>12</sup> was carefully removed from BioGRID prior to the current analysis. For each human protein interaction in BioGRID, a 2-hop neighborhood network was computed. Thus, for two pairs of interacting proteins A-B and B-C, the 2-hop network for protein A would include C. Bait-prey and prey-prey interactions in the current study were compared to the 2-hop network for corresponding proteins in BioGRID and the protein complexes in CORUM.

### *Protein-protein interaction false discovery rates*

The mean False Discovery Rate (FDR) estimate was computed by generating a null distribution of prey-prey cosine similarities from a reshuffled bait-prey D-score matrix (random permutation of rows and columns independently). This was repeated 100 times to estimate the standard error of the mean FDR. The number of prey-prey cosine similarity values greater than a given Prey-Prey Similarity (PPS) threshold was used to estimate the number of false prey-prey interactions. The FDR estimate is then the ratio of the number of false prey-prey interactions to the number of prey-prey interactions discovered. This strategy is similar to the ones used to estimate the FDR of peptide identifications from Peptide Spectrum Match with the help of decoy databases (reviewed in<sup>47</sup>), or to the re-sampling based methods

226 used by Lavalley-Adam et al.<sup>48</sup> and Dazard et al.<sup>49</sup> to assess the FDR of Protein-Protein Interactions  
 227 (PPI) from Affinity Purification Mass Spectrometry data with the help of a control set of experiments.  
 228 A controlled FDR of 7.58% +/-0.02% was achieved for the chosen PPS threshold of  $PPS \geq 0.75$  (see  
 229 Results section). Comparison plots of distributions of prey-prey cosine similarities in the original  
 230 matrix and in the null matrix under the re-sampling scheme are shown in Supplementary Figure S5. We  
 231 also looked at the FDR profile within a range of PPS, and found that the FDR was quite stable around  
 232 6-8% within the PPS range [0.5,0.75] and decreasing (as expected) within the PPS range [0.75,0.99].

### 233 *Semantic Similarity*

234 The semantic similarity is a metric for computing the similarity of Gene Ontology (GO) terms, their  
 235 ancestors, or the descendants for two genes<sup>24</sup>. Semantic similarity uses information content (IC) as a  
 236 measure of the specificity of a term and is quantified as the negative log likelihood,

$$237 \quad \quad \quad IC = -\log p(c) \quad \quad \quad (Eq. 4)$$

239  
 240 where  $p(c)$  is the probability of occurrence of  $c$  in the GO structure. We used Resnik's max measure of  
 241 similarity between two terms as the IC of their most informative common ancestor (MICA):

$$242 \quad \quad \quad sim_{Res}(c_1, c_2) = IC(c_{MICA}) \quad \quad \quad (Eq. 5)$$

244  
 245 This measure is effective in determining the information shared by the two terms<sup>25</sup>. Analysis was  
 246 performed using the *GOSim* R package<sup>26</sup>.

### 247 *Protein clustering and visualization*

248 Assembly of protein clusters from the prey-prey similarity matrix was performed using Pearson

249 correlation and hierarchical clustering (using TIBCO Spotfire software). Enrichment analyses of  
250 clustered sets of proteins were performed using FuncAssociate 2<sup>27</sup>. FuncAssociate performs a Fisher's  
251 Exact Test analysis to identify GO terms and the results are corrected for multiple hypotheses via  
252 empirical re-sampling, and adjusted p-values computed for significance. Cytoscape (version 2.8) was  
253 used for biological network visualization<sup>28</sup>.

## 254 **Results and Discussion**

### 255 *Initial data processing and analysis*

256 An overview of the data analysis workflow is shown in Figure 1. The workflow incorporates steps to  
257 address key issues in AP-MS datasets including the occurrence of non-specific prey proteins and  
258 missing data due to under-sampling of protein complexes. Scores for both bait-prey and prey-prey  
259 protein associations are computed and used to construct protein networks.  
260 Abundance and specificity of prey protein peptides in AP-MS experiments are important features that  
261 allow the confidence of prey proteins to be assessed. The D-score<sup>7</sup> was used as the primary metric for  
262 each bait-prey observation, since it combines measures of abundance and specificity into a single score  
263 and gives weight to well-replicated and less frequent prey proteins.

264 A principal challenge in interpreting AP-MS data is the presence of non-specific prey proteins.  
265 Typically, these are identified by their frequency across the dataset or as proteins present in control  
266 experiments. Our previous study<sup>12</sup> analyzed ~200 empty vector control AP-MS experiments that were  
267 then used to define non-specific prey proteins. We set a defined frequency threshold and used this  
268 threshold to filter out frequently occurring prey proteins<sup>12</sup>. In this analysis, we first observed that D-  
269 score values are typically very low for highly frequent proteins (since the D-score negatively weights  
270 frequency). Highly frequent proteins in the initial dataset such as PRMT5 (present in 94% of the 995  
271 AP-MS experiments in the dataset) have low median D-scores (the median D-score for PRMT5 is 1.8  
12



272 and 95% of baits for which PRMT5 is identified have D-score  $< 10$ ). The top 100 most frequent prey  
273 proteins in the dataset have a median D-score of  $\sim 3.0$  (Supplementary Figure S1). We used this value as  
274 a threshold to remove any bait-prey interaction with D-score  $\leq 3$ . Thus, 13688 bait-prey interactions  
275 with D-score  $\leq 3$  were removed from the dataset. The final matrix of D-scores used for subsequent  
276 analyses has dimensions of 384 baits x 5269 prey proteins, with  $\sim 1.25\%$  (25344) non-zero values  
277 (Supplementary Table S1).

278 To better identify relations between associated proteins in the sparse bait-prey matrix (Figure 1A), we  
279 used singular value decomposition (SVD), to project the bait-prey matrix into a lower dimensional  
280 space, and assign eigenvalues to bait-prey pairs (Figure 1B). Application of SVD to the bait-prey  
281 matrix increases the overall similarity of prey protein vectors within the matrix as shown in  
282 Supplementary Figure S2. SVD in the form of latent semantic analysis was previously applied to detect  
283 analytical trends in the HUPO Plasma Proteome Project (HUPO PPP)<sup>29</sup>. Our goal in applying SVD to  
284 the bait-prey matrix is to increase the probability of detecting biologically relevant associations  
285 between proteins in the matrix. The extent to which the matrix dimension is reduced is determined by  
286 the  $k$ -value in SVD (Figure 1B). A high  $k$ -value corresponds to a small reduction of matrix dimensions  
287 with possible retention of too much noise, whereas a small  $k$ -value may retain too little information  
288 from the original matrix. We estimated an appropriate  $k$ -value by plotting the singular value at the  $k^{\text{th}}$   
289 rank vs  $k$ -value (Supplementary Figure S3). A  $k$ -value of 150 was selected at the plateau in the graph so  
290 that the majority of significant information in the matrix is retained (singular values with rank  $< k$ ). To  
291 identify relationships between prey proteins in the bait-prey matrix we then generated a prey-prey  
292 similarity matrix (Figure 1C) by computing the similarity between each pair of preys in the matrix  
293 (cosine similarity). Although a similar methodology could also be applied to the bait-bait matrix, we  
294 focused on the prey-prey matrix (5269 x 5269 proteins), because it is much larger than the bait-bait  
295 matrix (384 x 384 proteins), and we reasoned that it would represent a richer resource for discovery of

13

novel interactions and complexes. To identify similarities between prey proteins, we computed the cosine distance between each pair of prey proteins for all prey proteins occurring with 2 or more baits (2242 preys) (Figure 1C).

*Benchmarking bait-prey and prey-prey interactions*

To gauge the biological relevance of the bait-prey and prey-prey interactions and their associated scores, we compared protein-protein pairs in our dataset to known datasets and measured the degree to which protein pairs shared functional annotations. Since the bait-prey and prey-prey interactions identified using our framework represent putative physical associations between proteins, we used two complementary sources of protein interaction data as our benchmark sources. CORUM<sup>22</sup> is a map of mammalian protein complexes curated from individual studies, whereas BioGRID<sup>23</sup> is a repository of physical protein interactions including data from high-throughput interaction studies. CORUM groups proteins according to protein complexes, whereas BioGRID represents protein interactions as protein-protein pairs. To benchmark versus BioGRID, we first computed the 2-hop neighborhood of each protein-protein pair in BioGRID. Since our data is derived from AP-MS experiments, representing complexes of physically associated proteins that may or may not directly interact, we compared our data to these neighborhoods of associated proteins rather than binary protein-protein pairs in BioGRID to ensure a representative comparison. For each bait-prey or prey-prey protein pair, we determined whether they co-occurred in a CORUM complex or within a BioGRID 2-hop network. To compare bait-prey and prey-prey interactions to known interactions, we calculated log likelihood scores for the relative enrichment of BioGRID and CORUM known interactions in our data, as previously formulated<sup>30</sup>. Figure 2 shows the log likelihood of known interactions for bait-prey D-scores (Figure 2A) and prey-prey scores (Figure 2B). Although the absolute numbers of interactions overlapping between our data and BioGRID are higher than the overlap between our data and

319 CORUM, CORUM interactions represented in our bait-prey or prey-prey datasets have significantly  
320 higher D-score or prey-prey scores respectively. Ranked by prey-prey score, the 25<sup>th</sup> percentile  
321 CORUM prey-prey score is 0.44 whereas the 25<sup>th</sup> percentile BioGRID prey-prey score is 0.24. This is  
322 not unexpected, since CORUM is derived from manually curated protein complexes whereas BioGRID  
323 includes high-throughput protein interaction studies. This also demonstrates the value of comparisons  
324 of protein interaction data to multiple sources; CORUM provides higher specificity with low sensitivity  
325 whereas BioGRID improves the sensitivity at the cost of lower specificity. Figure 2 also shows that the  
326 relative enrichment of BioGRID or CORUM interactions decreases as bait-prey D-score or prey-prey  
327 score decrease, showing that both scores provide some discrimination of true positive interactions from  
328 false. In the case of the bait-prey interactions, BioGRID and CORUM interactions decrease very  
329 sharply below D-score ~ 20 (corresponding to ~ 95<sup>th</sup> percentile of all of the D-scores, as determined by  
330 Sowa et al.<sup>7</sup>). In the case of the prey-prey interactions, BioGRID and CORUM interactions decrease  
331 more evenly across the range of prey-prey scores. The large predicted size of the human protein  
332 interactome<sup>30-31</sup>, the incompleteness of 'known' human protein interactions, as well as noise and  
333 context-sensitivity mean that intersections between experimental protein interaction datasets and  
334 known interactions tend to be small. For example, although BioGRID is a comprehensive source of  
335 available protein interaction data (~400,000 total interactions; ~38,000 human interactions)<sup>23</sup>, 35% of  
336 the bait proteins used in our original AP-MS study<sup>12</sup> have 1 or fewer interacting protein in BioGRID  
337 (27% have no interactions at all). The overlap between bait-prey or prey-prey interactions and known  
338 interactions in the BioGRID set or in CORUM was 26.9% of bait-prey interactions (D-score>20) and  
339 4.6% of prey-prey interactions (prey-prey score>0.75).

340 *Biological coherence of bait-prey and prey-prey associations*

341 To ascertain whether bait-prey and prey-prey protein pairs represent associations between proteins with

related functions, and to benchmark the bait-prey and prey-prey scores, we analyzed functional annotations of associated proteins using the Gene Ontology (GO)<sup>32</sup>. We first observed that computing co-annotation of GO terms for protein-protein pairs has low sensitivity, since many proteins, although biologically related may not be assigned the same term. We therefore used semantic similarity (SS) of GO terms which has proven to be robust measures of biological similarity for pairs or sets of genes<sup>24</sup>. Although there are multiple implementations of semantic similarity, here we use the Resnik max method<sup>25</sup>, since it was previously shown to perform best in a comparative analysis of semantic similarity metrics using large-scale protein-protein interactions<sup>33</sup>. Semantic similarity vastly increases the sensitivity of analyzing co-annotations over simple analysis of co-annotated proteins pairs (Supplementary Tables S2 and S3).

For both bait-prey and prey-prey associations, we reasoned that gene-ontology semantic similarity scores should be higher for bait-prey or prey-prey pairs with higher D-scores and prey-prey scores respectively. GO semantic similarity scores were computed for each bait-prey or prey-prey protein pair (Supplementary Tables S2 and S3) and analyzed as follows. Bait-prey protein pairs were binned according to D-score: high ( $D\text{-score} > 100$ ), medium ( $100.00 < D\text{-score} < 20.00$ ), low ( $20.0 < D\text{-Score} < 5.00$ ), very low ( $D\text{-score} < 5.00$ ) and a randomly selected set ( $n=1000$ ) of protein-protein pairs. Distributions of semantic similarity scores for bait-prey pairs are shown in Figure 3A for the molecular function GO ontology. The high and medium bins were found to be significantly higher than the random set (Wilcoxon test p-values :  $1.5E-08$ ,  $3.3E-06$ ,  $0.59$  for high, medium and low respectively). Similar trends were observed with the biological process and cellular compartment GO ontologies (Supplementary Figure S4). We used these analyses, along with the data shown in Figure 2A to calibrate the D-score threshold and therefore focused subsequent analyses on bait-prey pairs with D-score  $\geq 20$ .

Semantic similarity scores for all prey-prey associations were also computed and binned according to

366 prey-prey scores (Supplementary Table S3). Although not strictly monotonic, the log likelihood of  
367 known interaction enrichment with respect to prey-prey score broadly increases as prey-prey score  
368 increases (Figure 2B). We therefore grouped the prey-preys into four bins spanning the prey-prey  
369 score(PPS) range high ( $PPS \geq 0.75$ ), medium ( $0.75 < PPS < 0.50$ ), low ( $0.50 < PPS < 0.25$ ) and very low  
370 ( $PSS < 0.24$ ) and compared with semantic similarity of GO molecular function as shown in Figure 3B.  
371 A random set of 10000 prey-prey scores were generated to compute the significance test. Notably, the  
372 high scoring bin (prey-prey score  $\geq 0.75$ ) enriches for interactions with higher semantic similarity in  
373 all 3 GO ontologies (Figure S3), and the difference between the high bin and the random set was  
374 statistically significant (Wilcoxon Test p-values:  $2.1E-12$ ,  $0.54$  for high and medium respectively).  
375 These results show that both the D-score and prey-prey scores can be used to define sets of bait-prey or  
376 prey-prey pairs that are enriched for interactions with higher biological coherence, based upon their GO  
377 annotations. These analyses also provide a guide for selected subsets of interactions for further  
378 analysis, and as such we used bait-preys with  $D\text{-score} \geq 20$  and prey-preys with score  $\geq 0.75$  for  
379 building network models of protein complexes. Approximately 5.6% (1900) of bait-prey interactions  
380 and 7.7% (23,000) of prey-prey interactions meet these criteria, and were used in subsequent analyses.

### 381 *Identification of protein complexes*

382 Since AP-MS experiments identify co-complexed proteins, rather than binary pairs of interacting  
383 proteins, we organized the data into more meaningful biological groupings by clustering sets of  
384 proteins with significant prey-prey scores. In addition, since the sets of high-scoring bait-prey and prey-  
385 prey interactions are large, and likely contain significant numbers of false positives, clustering provides  
386 a means of focusing on higher-likelihood associations of proteins. The prey-prey similarity matrix was  
387 hierarchically clustered as shown in Figure 4. We identified 107 clusters comprised of a total of 754  
388 proteins (each cluster was required to have at least 3 proteins and all prey-prey associations  $> 0.9$ ) as

389 shown on the diagonal of the prey-prey matrix in Figure 4. For each of the 107 protein clusters, we  
390 identified significantly enriched GO categories, and for 43 of the 107 protein clusters, one or more  
391 significant ( $p < 0.05$ ) GO annotation terms were identified (Supplementary Table S4).

#### 392 *Network models incorporating bait-prey and prey-prey interactions*

393 As demonstrated in the previous sections, the global biological coherence of prey-prey interactions is  
394 similar to that of bait-prey interactions. Since matrix models of AP-MS data, in which all pairwise  
395 interactions are assumed to occur, are prone to false positives<sup>16</sup> we sought to selectively combine high-  
396 scoring bait-prey and prey-prey interactions into integrated network models. Protein clusters identified  
397 in Figure 4 were used as network seeds and extended by addition of selected high scoring bait-prey (D-  
398 score  $> 20$ ) and prey-prey (PPS  $> 0.75$ ) interactions. Four constructed network models, corresponding to  
399 Eukaryotic Initiation Factor (EIF) complexes, G-protein signaling and regulation, chromatin assembly  
400 factor complex and nucleosome regulation, and the proteasome are shown in Figure 5 (A-D  
401 respectively) and were selected to illustrate the potential of combining bait-prey and prey-prey  
402 interactions for delineation of network topology, and identification of new protein complex components  
403 and interactions.

404 Figure 5A shows a network model constructed by integrating bait-prey and prey-prey interactions  
405 corresponding to Eukaryotic Initiation Factor (EIF) complexes. Six bait proteins: EIF1B (also known as  
406 GC20), EIF2B1, EIF3S10, EIF4A2, EIF4A1, EIF4EBP1 and their associated prey proteins were  
407 integrated. Of particular interest is the separation of the cliques corresponding to EIF1/2/3 components  
408 (red shaded nodes) and EIF4 components (green nodes). A large number of prey proteins (D-score  $> 20$ )  
409 were found for EIF1B and EIF2B1 baits, including many ribosomal proteins, in line with known  
410 associations between these components and ribosomes<sup>34</sup>. In contrast, EIF4A2 and EIF4A1 baits had  
411 relatively few high-scoring prey proteins as shown in Figure 5A. Of particular note, the LSM14A

protein was identified uniquely with EIF4A2 bait and with high prey-prey interactions with EIF4G1 and PDCD4. LSM14A is a component of P-bodies, cytoplasmic foci that govern mRNA storage and degradation<sup>35</sup>. In line with previous findings, EIF4 components are present in P-bodies whilst other EIF components are conspicuously absent<sup>35</sup>. PDCD4 is uniquely associated with EIF4A2 and EIF4A1 baits in our data, and is known to interact directly with EIF4A2, EIF4A1 and EIF4G1<sup>36</sup>.

Whilst EIF complexes are well represented in our dataset, we also analyzed protein complexes with sparser coverage such as the networks in Figure 5B and Figure 5C. Figure 5B shows a network of proteins involved in G-protein signaling. Members of the Rho family of GTPases (RHOB, RHOC) and the Ras super-family of small GTP-binding proteins (RAC1, RAC2) are associated through their common bait, ARHGDIA. All of these proteins function in vesicular transport and endosomal signaling. An additional high-scoring interaction was observed with LYST, the lysosomal trafficking regulator. Endosomes may ultimately fuse with lysosomes for degradation of constituent molecules. In addition, LYST is associated with a rare lysosomal disorder, Chediak-Higashi syndrome, and thus the association with endosome signaling may potentially shed additional light on the disease mechanisms.

We also observed a cluster of prey-prey interactions corresponding to the Chromatin Assembly Factor (CAF-1) complex, and expanded this cluster of proteins into the network shown in Figure 5C. Several protein complexes with known functions were identified. The CAF-1 complex and Polycomb Recessive Complex 2 (PRC2) have related functions in chromatin metabolism and share components such as RBBP4. The recently identified MMS22L-TONSL complex<sup>37</sup> that mediates recombination mediated repair of replication forks and is comprised of MMS22L and TONSL (a.k.a NFKBIL2) proteins was also identified. Other significant interactions between these proteins and the Tousled-like kinases (TLK1, TLK2) were also observed. TLK1 and TLK2 heterodimerize and are also involved in chromatin assembly<sup>38</sup>. Within this rich, overlapping network of complexes with related functions, we searched for other high-scoring prey-prey interactions that might be novel components. Two other

436 proteins with significant prey-prey interactions with proteins in the chromatin modification network  
437 were observed. SMG6, a protein functioning in telomere regulation and nonsense-mediated mRNA  
438 decay was observed with high-scoring prey-prey interactions with several proteins (TLK1, 0.98;  
439 JARID2, 0.99). SMG6 is conserved (across eukaryotes) and has been found to physically interact with  
440 telomerase<sup>39</sup>. A second protein, UBN2, was recently identified as an ortholog of a yeast protein  
441 involved found in senescence associated chromatin foci<sup>40</sup>. The yeast ortholog of UBN2 interacts with  
442 the yeast ortholog of ASF1A, and both proteins function to create transcriptionally silent  
443 heterochromatin<sup>41</sup>. Thus, integration of bait-prey and prey-prey interactions serves to identify proteins  
444 linked to known complexes and functions.

445 Finally, we constructed a network based on coverage of the proteasome complex in our data (Figure  
446 5D) . This network was constructed from four bait proteins and 16 prey proteins, highly enriched for  
447 components of the eukaryotic proteasome. Three of the baits used (PSMD6, PSMD10 AND PSMD13)  
448 function as 26S proteasome non-ATPase regulatory subunits of the proteasome and whereas one bait  
449 (PSMC4) is a 26S protease regulatory subunit. Although the structure of the eukaryotic proteasome is  
450 comparatively well understood, we note that the network model delineates sub components of the  
451 proteasome. For example, PSMC1, PSMC2 and PSMC5 cluster whilst PSMD subunits (PSMD8,  
452 PSMD12 and PSMD14) cluster separately. Although the coverage of protein-protein interactions is  
453 relatively sparse, we note that the alpha-type proteasome subunits, PSMA5 and PSMA6 cluster  
454 exclusively with PSMD subunits and not PSMC subunits, possibly indicating intermediate forms of the  
455 proteasome comprised of specific sets of components. In summary, integrated network models as  
456 shown by these examples have the potential to yield novel components of protein complexes as well as  
457 further delineating the topology and sub-structure of networks.

458 We next compared the biological coherence of integrated (bait-prey and prey-prey) and matrix models  
459 of equivalent complexes. Matrix models for selected baits consisted of all pairwise protein interactions



(bait-prey and prey-prey), whereas integrated models consisted of all bait-prey interactions with D-score  $>20$  and all prey-prey interactions with score  $>0.75$ . Semantic similarity distributions were used to compare the integrated and matrix models of several protein complexes as shown in Figure 6. In addition to the Eukaryotic Initiation Factor (EIF) and Proteasomal complexes, we also analyzed data corresponding to four single baits (CTNNBIP1, VHL, REA and WDR8), the four most well replicated baits in the original study. For all four single baits and in the EIF complex, semantic similarity scores were significantly higher for integrated network models than for the matrix models. In only 1 case, the proteasome, the integrated and matrix model showed no significant difference between the semantic similarity distributions, suggesting that the biological coherence of integrated and matrix models of our proteasomal data is similar. This may be explained by the fact that most of the prey proteins identified by proteasome baits in our study are already known components of the proteasome, and so selection of prey-prey interactions with high PPS score for incorporation in the network model does not improve biological coherence over assuming that all preys interact with all other preys. In summary, these data show that selective integration of high-scoring bait-prey and prey-prey interactions can be used to generate protein network models with high biological coherence that can reveal new connections between proteins as well as new components of protein complexes.

## Conclusions

We present a data-driven framework for analysis of large-scale interaction proteomics data that uses integrated computational techniques to derive additional value from these datasets in terms of novel protein-protein interactions. We used this framework to analyze a unique, systematically generated human AP-MS dataset, in which complexes were determined for 384 disease-relevant bait proteins<sup>12</sup>. Specifically, our analysis integrates associations between the baits and their identified proteins (bait-prey), and associations that we detect amongst prey proteins (prey-prey), by analysis of the complete

483 data matrix. Analyzed globally, high-scoring bait-prey and prey-prey are enriched for known  
484 interactions and interactions between proteins with related biological function. In addition, by  
485 integrating prey-prey and bait-prey interactions into network models, we increase the biological  
486 coherence of those networks. We also show that integrated networks of bait-prey and prey-prey  
487 interactions provide a basis for delineation of network topology and identification of new protein  
488 complex members.

489 A major motivation for our work is to develop methods that enable novel protein-protein associations  
490 to be derived from large-scale datasets. Since mapping the complete human protein interactome  
491 experimentally is such an enormous undertaking, approaches that can be used to identify novel  
492 interactions from existing data, will continue to play a role in extending and refining the known human  
493 protein interactome. In addition, as previously observed<sup>22</sup>, despite the rapid accumulation of large  
494 volumes of proteomics data, there has been surprisingly little re-analysis and evaluation of most  
495 proteomics datasets. The exceptions to this include unique datasets such as the yeast interaction  
496 proteomics datasets<sup>2-5</sup> that have fueled much of the development of scoring algorithms as well as  
497 analysis of interaction networks.

498 Other studies with similar intention to the current work have primarily focused on the relatively  
499 complete yeast AP-MS datasets, or on AP-MS datasets that although not complete, are focused on  
500 specific complexes or sets of complexes<sup>6,20</sup>. In addition, most methods that have been proposed for AP-  
501 MS data analysis are not appropriate to all datasets, for reasons of data size or completeness. The  
502 heterogeneity of current AP-MS datasets in terms of biological and analytical methodology (epitope  
503 tag, expression construct, organism, cell-type etc) along with the challenges associated with the data  
504 itself (missing data, noise, contaminant proteins) have hindered the development of widely applicable  
505 analysis tools. Exceptions to this include the SAINT algorithm where the explicit goal is to develop a  
506 widely applicable AP-MS analysis method<sup>13</sup>.

With these challenges in mind, we have constructed a computational framework for analysis of large-scale human AP-MS data. Noteworthy features of our framework include representation of protein observations using the D-score that takes into account spectral counts, overall frequency, and replication for each protein observation. Although other studies with similar intent have made use of present/absent calls, particularly in analysis of the yeast AP-MS datasets, quantitative values derived via label-free analysis provide an additional dimension for ranking and discrimination of true positive interactions<sup>20</sup>. In the latter study, Choi *et al* used nested bi-clustering to group sets of baits or preys with similar quantitative spectral count levels. In our study, the dataset is more heterogeneous both in terms of the actual baits used as well as the number of replicate AP-MS experiments per bait. For this reason, rather than use the spectral counts directly for each protein, we used the D-score as the starting point for the analysis, so that frequency and number of replicates as well as spectral counts are taken into account. Second, we address the issue of sparseness of the primary bait-prey matrix by first transforming the matrix of bait-prey scores through singular value decomposition. Singular value decomposition in the form of latent semantic analysis has previously been used to large-scale expression proteomics albeit with a different goal<sup>29</sup>. In the latter study, SVD was used to analyze heterogeneous plasma proteomics datasets acquired using different technologies in different laboratories. The method used by Klie *et al* represented protein observations as binary measures, presumably so that data from different instruments and different laboratories could be appropriately integrated. Regardless of whether quantitative values are used however, tools such as SVD, that address the sparseness and missing data challenges of proteomics datasets are essential. SVD was also previously applied to the assembly of protein interaction networks from more focused human AP-MS data, although in this case, SVD was applied to the problem of identifying clusters of proteins<sup>6</sup>. To mitigate the problems of false positives in studying the large volume of prey-prey associations, we use semantic similarity measures of annotation to first calibrate our protein-protein scores (bait-prey

and prey-prey) and then to select small high-scoring subsets of protein-protein interactions to study. A principal challenge of data-driven inference of human protein-protein interactions remains the lack of 'gold-standard' protein interactions and complexes. Although curated annotations of mammalian protein complexes exist<sup>22</sup>, these represent only a fraction of the total interactome. In the absence of gold-standard datasets, the semantic similarity measures provide a means of benchmarking protein-protein scores as shown here. Our method increases the amount of information in terms of protein-protein associations that may be gleaned from large-scale AP-MS studies. However, increasing the density of coverage (in terms of bait proteins used) will be the method of choice for truly defining protein complexes in the cell. Studies that iterate through a network, testing all proteins as baits in AP-MS experiments may provide the most detailed representations of protein complexes<sup>8</sup>. The complexity of mapping the protein interactome, in terms of distinguishing complexes that share components was highlighted in an AP-MS study of chromatin remodeling complexes<sup>6</sup>. This study reiterates the point that although computational and statistical analyses may allow for the prediction of protein-protein interactions and network topology, a detailed map of overlapping protein complexes may ultimately only be achieved through high density experimental analyses.

Integration of other data may further help refine of protein network topology. Although protein quantifications are only loosely correlated with protein-protein interactions, integration of protein abundance values may be used to refine prediction of protein interaction. The recently developed PaxDb<sup>45</sup> is a cross-species database of protein quantifications, thus providing an independent source of proteome-wide abundance. In PaxDb, protein-protein interactions are used as a metric of consistency for quantitative proteomics studies, with the rationale that interacting proteins tend to have more similar levels of expression than randomly selected or non-interacting proteins<sup>45</sup>. Proteome-wide protein abundance values may also prove useful for normalizing the abundance values of proteins identified in AP-MS experiments. For example, a recent study used PaxDb values to account for the differing

cellular abundance of proteins identified in AP-MS analysis of chromatin remodeling complexes<sup>46</sup>. Future work might therefore utilize the approach that we have described in conjunction with other proteome-wide information for further refinement of protein networks and protein interaction discovery.

## **Supplementary Material**

**Supplementary Figure S1** The density plot of top 100 frequent control proteins D-scores in the bait-prey dataset. The top 100 frequent control proteins were selected from the 200 control experiments (without bait), and the median D-scores of these 100 frequent control proteins in this study was close to 3.

**Supplementary Figure S2** The two-dimensional plot of prey-prey similarities scores using SVD in the bait-prey matrix versus without using SVD. The application of SVD to the bait-prey matrix increases the overall similarity of prey protein vectors within the matrix.

**Supplementary Figure S3** Selection of  $K$ -value for singular value decomposition of bait prey score matrix. Singular values vs.  $K$  values was used to identify optimum  $K$  value=150. The value of  $K$  determines the degree of reduction, a high  $K$  value corresponds to small reduction (minimal filtering of noise), while a small  $K$  value correspond to strong reduction (too little information retained)

**Supplementary Figure S4** GO molecular function (MF), biological process (BP) and cellular compartment (CC) semantic similarities comparison with bait-prey D-scores and prey-prey similarity scores. A. Box plots represent distributions of semantic similarity scores for bait-prey protein pairs binned according to D-scores (4 bins) for gene ontologies. Bin 1 ( $D\text{-score} > 100$ ) and bin 2 ( $100 < D\text{-score} > 20$ ) have statistically significantly higher semantic similarity than random set (R) of D-scores in

579 MF, BP and CC (Wilcoxon Test P values are 1.5E-8, 7.9E-6 and 1.7E-7 respectively in bin 1). B. Box  
580 plots represent distributions of semantic similarity scores for prey-prey protein pairs binned according  
581 to prey-prey scores (4 bins) for gene ontologies. Bin 1 prey-prey scores ( $>0.75$ ) bin has significantly  
582 higher semantic similarity than random set (R) of prey-prey similarity scores in MF, BP and CC (  
583 Wilcoxon Test P values are 2.1E-12, 7.2E-12 and 5.0E-4 respectively) . Median of each distribution is  
584 represented by horizontal bar in each box plot.

585

586 **Supplementary Figure S5** Exploratory Data Analysis plot of the distributions of prey-prey cosine  
587 similarity values. Histograms represent the distribution of prey-prey cosine similarity values in the  
588 original matrix (left) and in the null matrix under the resampling scheme (middle). The quantile-  
589 quantile plot of prey-prey cosine similarities (right) shows the strong deviation of the two distributions  
590 in the two situations.

591 **Supplementary Table S1** Bait-prey D-score table. Bait protein purifications name, prey protein  
592 identified, and D-score.

593 **Supplementary Table S2** Comparison Bait-Prey (*A* and *B*) D-score with Molecular Function (MF),  
594 Biological Process(BP) and Cellular Compartment(CC) GO terms semantic similarity. The column  
595 names of this table are: Gene symbol *A*, Entrez gene ID *A*, Gene symbol *B*, Entrez gene ID *B*, Bait *A*-  
596 Prey *B* D-score, Rank (based on Bait-Prey D-score, where 1=High (D-score $>100$ ); 2=medium (100<D-  
597 score $>20$ ); 3=Low (20<D-score $>5$ ); 4=very Low (D-score $<5$ )), Bait *A*- Prey *B* MF semantic similarity  
598 score, Bait *A*- Prey *B* BP semantic similarity score, Bait *A*- Prey *B* CC semantic similarity score.

599 **Supplementary Table S3** Comparison Prey-Prey (*A* and *B*) similarity score with Molecular Function  
600 (MF), Biological Process(BP) and Cellular Compartment(CC) GO semantic similarity. The column  
601 names of this table are: Gene symbol *A*, Entrez gene ID *A*, Gene symbol *B*, Entrez gene ID *B*, Prey *A*-  
602 Prey *B* similarity score, Rank (based on Prey-Prey similarity score(PPSS), where 1=High(PPSS $>0.75$ );  
26

2=medium( $0.75 < \text{PPSS} < 0.5$ ); 3=Low ( $0.5 < \text{PPSS} < 0.25$ ); 4=very low ( $\text{PPSS} < 0.25$ )), MF semantic similarity score, BP semantic similarity score, CC semantic similarity score.

**Supplementary Table S4** Protein complexes identified. Modules name describing its components along with common bait proteins purification names, enriched GO terms, GO IDs , significant P-values.

**Acknowledgments**

We thank Joseph Abraham and Parminder Kaur for productive scientific discussions. R. M. E. acknowledges funds from the Cleveland Foundation used in part to fund this study.

**Supporting Information Available:** This material is available free of charge via the Internet at <http://pubs.acs.org>.

**References**

- Saha S, Kaur P, Ewing RM. The bait compatibility index: computational bait selection for interaction proteomics experiments. *J Proteome Res.* 2010 Oct 1;9(10):4972-81.
- Ho Y, Gruhler A, Heilbut A, Bader GD, Moore L, Adams SL, Millar A, Taylor P, Bennett K, Boutilier K, Yang L, Wolting C, Donaldson I, Schandorff S, Shewnarane J, Vo M, Taggart J, Goudreault M, Muskat B, Alfarano C, Dewar D, Lin Z, Michalickova K, Willems AR, Sassi H, Nielsen PA, Rasmussen KJ, Andersen JR, Johansen LE, Hansen LH, Jespersen H, Podtelejnikov A, Nielsen E, Crawford J, Poulsen V, Sørensen BD, Matthiesen J, Hendrickson RC, Gleeson F, Pawson T, Moran MF, Durocher D, Mann M, Hogue CW, Figeys D, Tyers M. Systematic identification of protein complexes in *Saccharomyces cerevisiae* by mass spectrometry. *Nature.* 2002 Jan 10;415(6868):180-3.
- GavinAC, Bösch M, Krause R, Grandi P, Marzioch M, Bauer A, Schultz J, Rick JM, Michon AM, Cruciat CM, Remor M, Höfert C, Schelder M, Brajenovic M, Ruffner H, Merino A, Klein K, Hudak M, Dickson D, Rudi T, Gnau V, Bauch A, Bastuck S, Huhse B, Leutwein C, Heurtier MA, Copley RR, Edelmann A, Querfurth E, Rybin V, Drewes G, Raida M, Bouwmeester T, Bork P, Seraphin B, Kuster B, Neubauer G, Superti-Furga G. Functional organization of the yeast proteome by systematic analysis of protein complexes. *Nature.* 2002 Jan 10;415(6868):141-7.
- Krogan NJ, Cagney G, Yu H, Zhong G, Guo X, Ignatchenko A, Li J, Pu S, Datta N, Tikuisis AP, Punna T, Peregrín-Alvarez JM, Shales M, Zhang X, Davey M, Robinson MD, Paccanaro A, Bray JE, Sheung A, Beattie B, Richards DP, Canadien V, Lalev A, Mena F, Wong P, Starostine A, Canete MM, Vlasblom J, Wu S, Orsi C, Collins SR, Chandran S, Haw R, Rilstone JJ, Gandi K, Thompson NJ, Musso G, St Onge P, Ghanny S, Lam MH, Butland G, Altaf-UI AM, Kanaya S, Shilatifard A, O'Shea E, Weissman JS, Ingles CJ, Hughes TR, Parkinson J, Gerstein M, Wodak SJ, Emili A, Greenblatt JF. Global landscape of protein complexes in the yeast *Saccharomyces cerevisiae*. *Nature.* 2006 Mar 30;440(7084):637-43
- Yu H, Braun P, Yildirim MA, Lemmens I, Venkatesan K, Sahalie J, Hirozane-Kishikawa T, Gebreab F, Li N, Simonis N, Hao T, Rual JF, Dricot A, Vazquez A, Murray RR, Simon C, Tardivo L, Tam S, Svrikapa N,

639 Fan C, de Smet AS, Motyl A, Hudson ME, Park J, Xin X, Cusick ME, Moore T, Boone C, Snyder M, Roth  
640 FP, Barabási AL, Tavernier J, Hill DE, Vidal M. High-quality binary protein interaction map of the yeast  
641 interactome network. *Science*. 2008 Oct 3;322(5898):104-10.

642 6. Sardiú ME, Cai Y, Jin J, Swanson SK, Conaway RC, Conaway JW, Florens L, Washburn MP.  
643 Probabilistic assembly of human protein interaction networks from label-free quantitative proteomics.  
644 *Proc Natl Acad Sci U S A*. 2008 Feb 5;105(5):1454-9.

645 7. Sowa ME, Bennett EJ, Gygi SP, Harper JW. Defining the human deubiquitinating enzyme interaction  
646 landscape. *Cell*. 2009 Jul 23;138(2):389-403.

647 8. Goudreault M, D'Ambrosio LM, Kean MJ, Mullin MJ, Larsen BG, Sanchez A, Chaudhry S, Chen GI,  
648 Sicheri F, Nesvizhskii AI, Aebersold R, Raught B, Gingras AC. A PP2A phosphatase high density  
649 interaction network identifies a novel striatin-interacting phosphatase and kinase complex linked to the  
650 cerebral cavernous malformation 3 (CCM3) protein. *Mol Cell Proteomics*. 2009 Jan;8(1):157-71.

651 9. Glatter T, Wepf A, Aebersold R, Gstaiger M. An integrated workflow for charting the human interaction  
652 proteome: insights into the PP2A system. *Mol Syst Biol*. 2009;5:237.

653 10. Gavin AC, Aloy P, Grandi P, Krause R, Boesche M, Marzioch M, Rau C, Jensen LJ, Bastuck S,  
654 Dimpelfeld B, Edelmann A, Heurtier MA, Hoffman V, Hoefert C, Klein K, Hudak M, Michon AM, Schelder  
655 M, Schirle M, Remor M, Rudi T, Hooper S, Bauer A, Bouwmeester T, Casari G, Drewes G, Neubauer G,  
656 Rick JM, Kuster B, Bork P, Russell RB, Superti-Furga G. Proteome survey reveals modularity of the  
657 yeast cell machinery. *Nature*. 2006 Mar 30;440(7084):631-6.

658 11. Collins SR, Kemmeren P, Zhao XC, Greenblatt JF, Spencer F, Holstege FC, Weissman JS, Krogan NJ.  
659 Toward a comprehensive atlas of the physical interactome of *Saccharomyces cerevisiae*. *Mol Cell*  
660 *Proteomics*. 2007 Mar;6(3):439-50.

661 12. Ewing RM, Chu P, Elisma F, Li H, Taylor P, Climie S, McBroom-Cerajewski L, Robinson MD, O'Connor L,  
662 Li M, Taylor R, Dharsee M, Ho Y, Heilbut A, Moore L, Zhang S, Ornatsky O, Bukhman YV, Ethier M,  
663 Sheng Y, Vasilescu J, Abu-Farha M, Lambert JP, Duwel HS, Stewart II, Kuehl B, Hogue K, Colwill K,  
664 Gladwish K, Muskat B, Kinach R, Adams SL, Moran MF, Morin GB, Topaloglou T, Figeys D. Large-scale  
665 mapping of human protein-protein interactions by mass spectrometry. *Mol Syst Biol*. 2007;3:89.

666 13. Choi H, Larsen B, Lin ZY, Breitkreutz A, Mellacheruvu D, Fermin D, Qin ZS, Tyers M, Gingras AC,  
667 Nesvizhskii AI. SAINT: probabilistic scoring of affinity purification-mass spectrometry data. *Nat Methods*.  
668 2011 Jan;8(1):70-3.

669 14. Lavallée-Adam M, Cloutier P, Coulombe B, Blanchette M. Modeling contaminants in AP-MS/MS  
670 experiments. *J Proteome Res*. 2011 Feb 4;10(2):886-95.

671 15. Skarra DV, Goudreault M, Choi H, Mullin M, Nesvizhskii AI, Gingras AC, Honkanen RE. Label-free  
672 quantitative proteomics and SAINT analysis enable interactome mapping for the human Ser/Thr protein  
673 phosphatase 5. *Proteomics*. 2011 Apr;11(8):1508-16. doi: 10.1002/pmic.201000770.

674 16. Bader GD, Hogue CW. Analyzing yeast protein-protein interaction data obtained from different sources.  
675 *Nat Biotechnol*. 2002 Oct;20(10):991-7.

676 17. Hart GT, Lee I, Marcotte ER. A high-accuracy consensus map of yeast protein complexes reveals  
677 modular nature of gene essentiality. *BMC Bioinformatics*. 2007 Jul 2;8:236.

678 18. Yu X, Ivanic J, Wallqvist A, Reifman J. A novel scoring approach for protein co-purification data reveals  
679 high interaction specificity. *PLoS Comput Biol*. 2009 Sep;5(9):e1000515.

680 19. Geva G, Sharan R. Identification of protein complexes from co-immunoprecipitation data. *Bioinformatics*.  
681 2011 Jan 1;27(1):111-7.

682 20. Choi H, Kim S, Gingras AC, Nesvizhskii AI. Analysis of protein complexes through model-based  
683 biclustering of label-free quantitative AP-MS data. *Mol Syst Biol*. 2010 Jun 22;6:385

684 21. Deerwester S, Dumais St, Furnas GW, Landauer TK and Harshman R. Indexing by Latent Semantic  
685 Analysis. *J. Am. Soc. Inf. Sci*. 1990, 41(6), 391-407.



22. Ruepp A, Waegelé B, Lechner M, Brauner B, Dunger-Kaltenbach I, Fobo G, Frishman G, Montrone C, Mewes HW. CORUM: the comprehensive resource of mammalian protein complexes—2009. *Nucleic Acids Res.* 2010 Jan;38(Database issue):D497-501.
23. Stark C, Breitkreutz BJ, Chatr-Aryamontri A, Boucher L, Oughtred R, Livstone MS, Nixon J, Van Auken K, Wang X, Shi X, Reguly T, Rust JM, Winter A, Dolinski K, Tyers M. The BioGRID Interaction Database: 2011 update. *Nucleic Acids Res.* 2011 Jan;39(Database issue):D698-704.
24. Pesquita C, Faria D, Falcão AO, Lord P, Couto FM. Semantic similarity in biomedical ontologies. *PLoS Comput Biol.* 2009 Jul;5(7):e1000443.
25. Resnik P (1995) Using information content to evaluate semantic similarity in a taxonomy. In: *Proc. Of the 14<sup>th</sup> International Jpint Conference on Artificial Intelligence.* Pp 448-453.
26. Fröhlich H, Speer N, Poustka A, Beissbarth T. GOSim--an R-package for computation of information theoretic GO similarities between terms and gene products. *BMC Bioinformatics.* 2007 May 22;8:166.
27. Berriz GF, Beaver JE, Cenik C, Tasan M, Roth FP. Next generation software for functional trend analysis. *Bioinformatics.* 2009 Nov 15;25(22):3043-4.
28. Smoot ME, Ono K, Ruscheinski J, Wang PL, Ideker T. Cytoscape 2.8: new features for data integration and network visualization. *Bioinformatics.* 2011 Feb 1;27(3):431-2.
29. Klie S, Martens L, Vizcaíno JA, Côté R, Jones P, Apweiler R, Hinneburg A, Hermjakob H. Analyzing large-scale proteomics projects with latent semantic indexing. *J Proteome Res.* 2008 Jan;7(1):182-91.
30. Ramani AK, Bunescu RC, Mooney RJ, Marcotte EM. Consolidating the set of known human protein-protein interactions in preparation for large-scale mapping of the human interactome. *Genome Biol.* 2005;6(5):R40.
31. Stumpf MP, Thorne T, de Silva E, Stewart R, An HJ, Lappe M, Wiuf C. Estimating the size of the human interactome. *Proc Natl Acad Sci U S A.* 2008 May 13;105(19):6959-64.
32. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight SS, Eppig JT, Harris MA, Hill DP, Issel-Tarver L, Kasarskis A, Lewis S, Matese JC, Richardson JE, Ringwald M, Rubin GM, Sherlock G. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet.* 2000 May;25(1):25-9.
33. Xu T, Du L, Zhou Y. Evaluation of GO-based functional similarity measures using *S. cerevisiae* protein interaction and expression profile data. *BMC Bioinformatics.* 2008 Nov 6;9:472.
34. Pestova TV, Kolupaeva VG, Lomakin IB, Pilipenko EV, Shatsky IN, Agol VI, Hellen CU. Molecular mechanisms of translation initiation in eukaryotes. *Proc Natl Acad Sci U S A.* 2001 Jun 19;98(13):7029-36.
35. Eulalio A, Behm-Ansmant I, Schweizer D, Izaurralde E. P-body formation is a consequence, not the cause, of RNA-mediated gene silencing. *Mol Cell Biol.* 2007 Jun;27(11):3970-81.
36. Yang HS, Jansen AP, Komar AA, Zheng X, Merrick WC, Costes S, Lockett SJ, Sonenberg N, Colburn NH. The transformation suppressor Pdc4 is a novel eukaryotic translation initiation factor 4A binding protein that inhibits translation. *Mol Cell Biol.* 2003 Jan;23(1):26-37.
37. Piwko W, Olma MH, Held M, Bianco JN, Pedrioli PG, Hofmann K, Pasero P, Gerlich DW, Peter M. RNAi-based screening identifies the Mms22L-Nfkbil2 complex as a novel regulator of DNA replication in human cells. *EMBO J.* 2010 Dec 15;29(24):4210-22.
38. Groth A, Lukas J, Nigg EA, Silljé HH, Wernstedt C, Bartek J, Hansen K. Human Tosl-like kinases are targeted by an ATM- and Chk1-dependent DNA damage checkpoint. *EMBO J.* 2003 Apr 1;22(7):1676-87.
39. Redon S, Reichenbach P, Lingner J. Protein RNA and protein protein interactions mediate association of human EST1A/SMG6 with telomerase. *Nucleic Acids Res.* 2007;35(20):7011-22.
40. Banumathy G, Somaiah N, Zhang R, Tang Y, Hoffmann J, Andrade M, Ceulemans H, Schultz D, Marmorstein R, Adams PD. Human UBN1 is an ortholog of yeast Hpc2p and has an essential role in the HIRA/ASF1a chromatin-remodeling pathway in senescent cells. *Mol Cell Biol.* 2009 Feb;29(3):758-70.

733 41. Zhang R, Poustovoitov MV, Ye X, Santos HA, Chen W, Daganzo SM, Erzberger JP, Serebriiskii IG,  
734 Canutescu AA, Dunbrack RL, Pehrson JR, Berger JM, Kaufman PD, Adams PD. Formation of  
735 MacroH2A-containing senescence-associated heterochromatin foci and senescence driven by ASF1a  
736 and HIRA. *Dev Cell*. 2005 Jan;8(1):19-30.

737 42. Xu, H., Freitas, M. A., MassMatrix: A database search program for rapid characterization of proteins and  
738 peptides from tandem mass spectrometry data. *Proteomics* **2009**, 9, (6), 1548–1555.

739 43. Nesvizhskii, A. I., Keller, A., Kolker, E. & Aebersold, R. A statistical model for identifying proteins by  
740 tandem mass spectrometry *Anal Chem*, **2003**, 75, 4646-58.

741 44. Nelson, E. K., Piehler, B., Eckels, J., Rauch, A., Bellew, M., Hussey, P., Ramsay, S., Nathe, C., Lum, K.,  
742 Krouse, K., Stearns, D., Connolly, B., Skillman, T. & Igra, M. LabKey Server: an open source platform for  
743 scientific data integration, analysis and collaboration. *BMC Bioinformatics*, **2011**, 12, 71

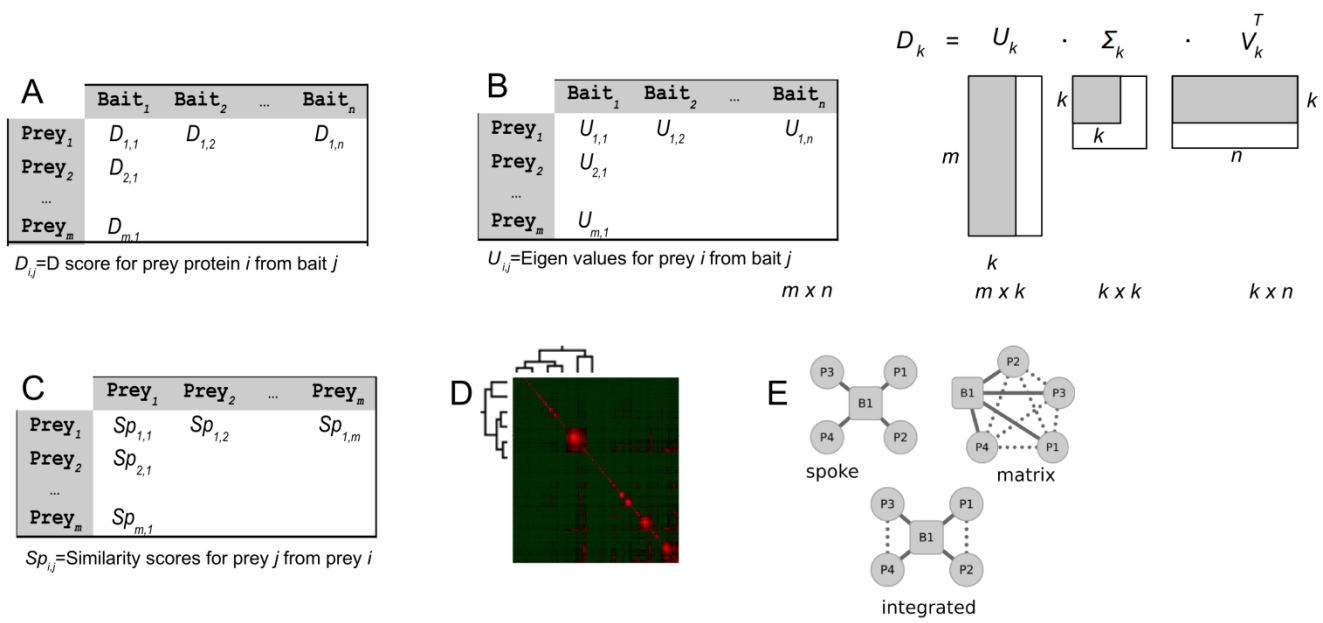
744 45. Wang, M.; Weiss, M.; Simonovic, M.; Haertinger, G.; Schrimpf, S. P.; Hengartner, M. O.; von Mering, C.  
745 *Molecular & cellular proteomics* **2012**. 10.1074/mcp.O111.014704

746 46. Tsai, Y.-C.; Greco, T. M.; Boonmee, A.; Miteva, Y.; Cristea, I. M. *Molecular & cellular proteomics* **2012**, 11,  
747 M111.015156.

748 47. Nesvizhskii, A. I., A survey of computational methods and error rate estimation procedures for peptide  
749 and protein identification in shotgun proteomics. *J Proteomics* 2010, 73(11), 2092-123.

750 48. Lavalley-Adam, M.; Cloutier, P.; Coulombe, B.; Blanchette, M., Modeling contaminants in AP-MS/MS  
751 experiments. *Journal of proteome research* 2011, 10(2), 886-95.

752 49. Dazard, J. E.; Saha, S.; Ewing, R. M., ROCS: A reproducibility index and confidence score for interaction  
753 proteomics. *BMC bioinformatics* 2012, 13(1), 128.



755

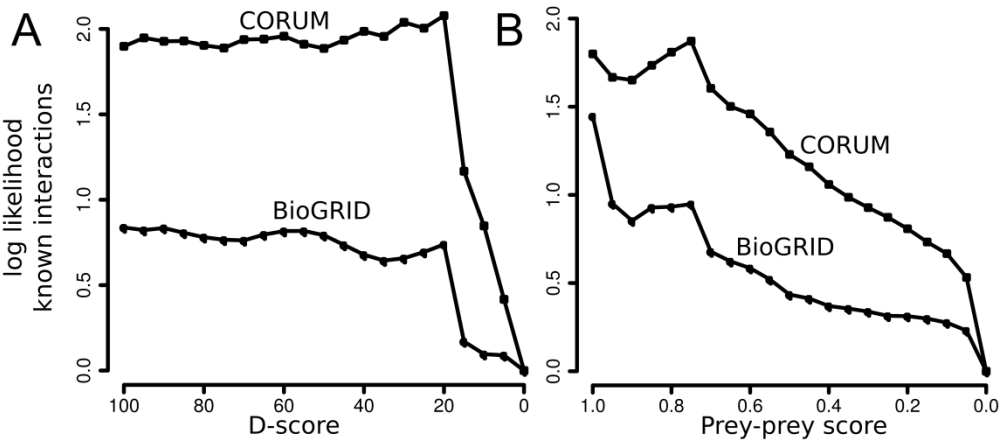
756

757 **Figure 1. Data analysis work-flow.** **A.** D-score<sup>7</sup> matrix with score for each pair of bait-prey proteins.  
758 **B.** D-score matrix approximation using singular value decomposition. **C.** Pairwise cosine similarity  
759 computed for each vector of prey scores. **D.** Topological overlap matrix created using hierarchical  
760 clustering to group prey proteins with similar prey-prey profiles. **E.** Spoke, matrix and integrated  
761 models for a hypothetical bait (B1) and 4 prey proteins (P1-P4). Integrated model incorporates selected  
762 edges from spoke and matrix models (solid lines lines represent bait-prey interactions, dotted lines  
763 represent prey-prey interactions). Figure 1B adapted with permission from (Klie et al, Journal of  
764 Proteome Research 7, 182-19). Copyright (2008) American Chemical Society.

765

766

767

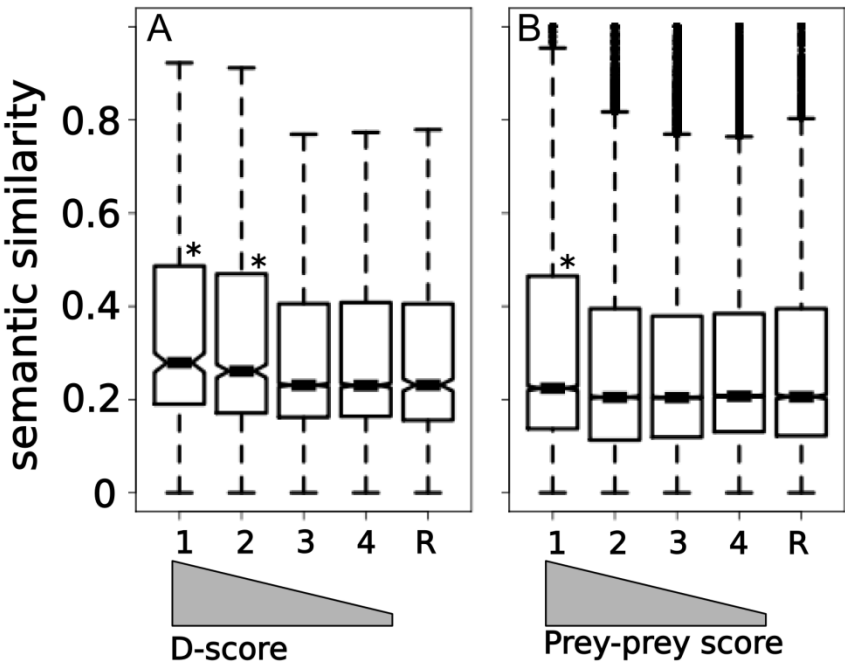


768

769 **Figure 2. Comparison with known interactions.** A. Bait-prey interactions compared to known  
770 interactions. Bait-prey protein pairs within shared 2-hop BioGRID network or present in same  
771 CORUM complex counted as 'known'. Log likelihood of relative enrichment of known vs. unknown  
772 interactions computed for each D-score threshold. B. As A, with log likelihood computed for each  
773 prey-prey protein pair and corresponding prey-prey interaction score threshold.

774

775



777

778

779 **Figure 3. GO semantic similarity distributions and protein interaction scores. A.** Box plots

780 represent distributions of semantic similarity scores for bait-prey protein pairs binned according to D-

781 scores (4 bins) for molecular function gene ontologies. Bins 1-4 represent D-score bins of D-score >

782 100, 100>D-score>20, 20>D-score>5 and D-score<5 respectively. Bins 1 and 2 have statistically

783 significantly (\*) higher semantic similarity than random set (R) (Wilcoxon Test P-values are 1.5E-08

784 and 3.3E-06). **B.** Box plots represent distributions of semantic similarity scores for prey-prey protein

785 pairs binned according to prey-prey scores (4 bins) for molecular function (MF) gene ontologies. Bins

786 1-4 represent Prey-prey scores of >0.75, 0.75>PPS>0.5, 0.5>PPS>0.25 and PPS<0.25 respectively. Bin

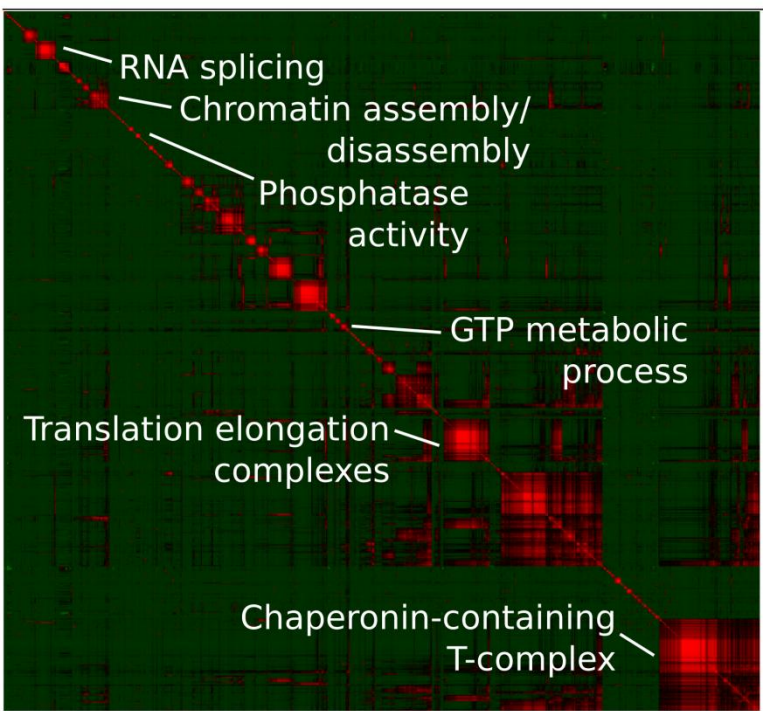
787 1 has significantly (\*) higher semantic similarity than random set (R) of prey-prey similarity scores

788 (Wilcoxon Test P value is 2.1E-12). Median of each distribution is represented by horizontal bar in each

789 box plot.

790

791



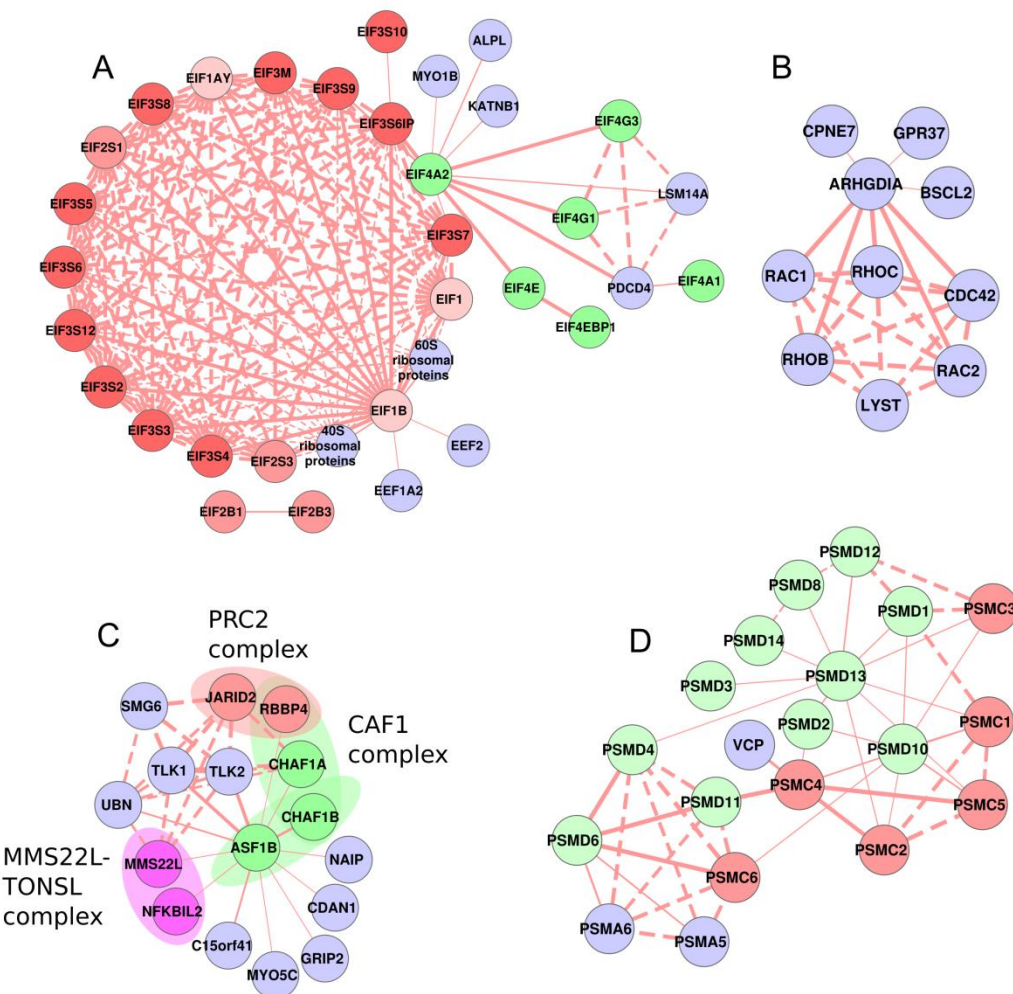
792

793

794 **Figure 4. Topological overlap matrix of prey-prey interactions.** The prey-prey similarity matrix  
795 (2242 X 2242) was used to create a topological overlap matrix using hierarchical clustering to group  
796 preys with similar prey-prey vectors. Red areas indicate prey-prey vector similarity and diagonal shows  
797 107 protein modules. Selected modules are labeled with significant GO terms.

798

799



801

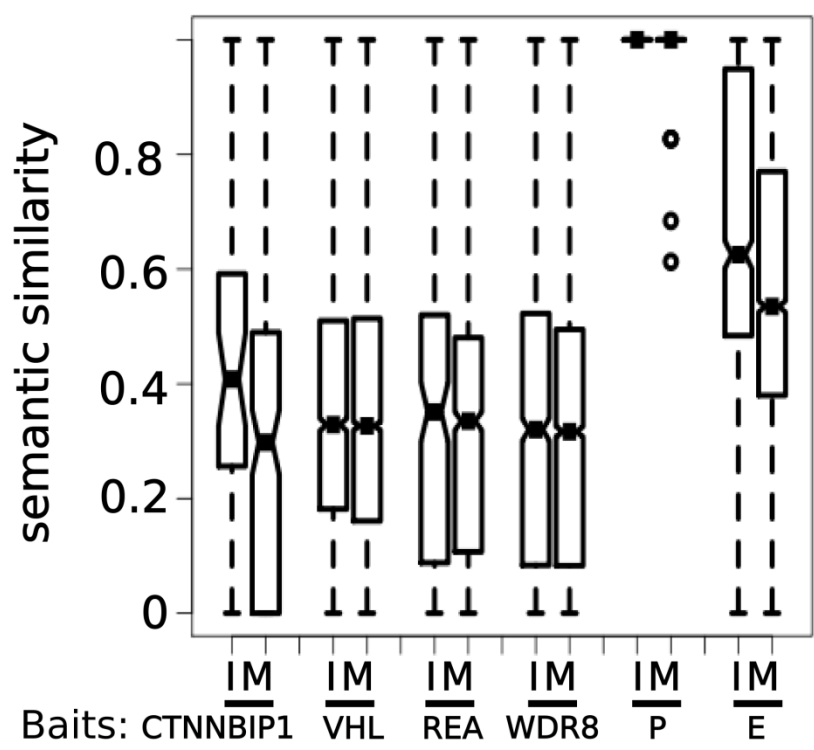
802

803

804 **Figure 5. Selected network models integrating high-scoring bait-prey and prey-prey edges. A.**  
805 Eukaryotic Initiation Factor (EIF) bait proteins (EIF1B, EIF3S10, EIF2B1, EIF4EBP1, EIF4A1 and  
806 EIF4A2) and their associated prey proteins (D-score>20) were combined with associated high-scoring  
807 prey-prey interactions. EIF1/2/3, EIF4 and EIF2 components are represented in red, green and gray  
808 nodes respectively. **B.** G-protein signaling complex. ARHGDIA bait and associated GTP binding  
809 proteins, RAC1 and RAC2; GTPase, RHOB and RHOC are identified preys with extensive inter-

810 connectivity (high prey-prey scores). **C.** Chromatin metabolism/ CAF1 complex, PRC2 complex and  
811 MMS22L-TONSL complex are shown in green, red and pink color respectively. **D.** Proteasome  
812 complex, four bait proteins corresponding to proteasome, PSMD6, PSMD10, PSMD13 and PSMC4  
813 and 16 prey proteins forming highly enriched eukaryotic proteasome complex. The gray nodes  
814 correspond to proteasome core complex subunit (PSMA5/6); green nodes to 26S proteasome non-  
815 ATPase regulatory subunit (PSMD\*); red nodes represents 26S protease regulatory subunits (PSMC\*).  
816 Solid edges indicate bait-prey interactions and dashed edges indicate prey-prey interactions, edge  
817 thickness indicates bait-prey or prey-prey scores appropriately (thicker edge/higher score).  
818  
819





834

835

836

837 **Figure 6. Biological coherence of integrated and matrix network models of protein complexes.**

838 Semantic similarity (GO biological process) distributions of protein-protein interactions in Integrated  
839 (I) and Matrix (M) models. Networks were constructed for data from 4 single baits (CTNNBIP1, VHL,  
840 REA and WDR8) and two models incorporating multiple baits, Proteasome complex (P) and  
841 Eukaryotic Initiation Factor (E) complexes. In 5 cases (CTNNBIP1, VHL, REA, WDR8, E), the  
842 Integrated model shows higher median semantic similarity than the Matrix model, and in 4 of these  
843 cases (CTNNBIP1,VHL,WDR8, E) the difference between Integrated and Matrix models is significant  
844 (Wilcoxon Test P-values<0.1).

845