

1 **Systematic protein complex profiling and differential analysis from co-** 2 **fractionation mass spectrometry data**

3
4 Andrea Fossati^{1,2,†}, Chen Li^{1,3,†,*}, Peter Sykacek⁴, Moritz Heusel⁵, Fabian Frommelt¹, Federico Uliana¹,
5 Mahmoud Hallal⁶, Isabell Bludau^{1,7}, Capraz Tümay Klemens⁴, Peng Xue^{1,8}, Anthony W. Purcell³,
6 Matthias Gstaiger^{1,*}, and Ruedi Aebersold^{1,9,*}

7
8 ¹Department of Biology, Institute of Molecular Systems Biology, ETH Zürich, Switzerland;
9 ²Department of Biology, Institute of Molecular Health Sciences, ETH Zürich, Switzerland;
10 ³Department of Biochemistry and Molecular Biology and Infection and Immunity Program,
11 Biomedicine Discovery Institute, Monash University, VIC 3800, Australia; ⁴Department of
12 Biotechnology, BOKU University, Vienna, Austria; ⁵Division of Infection Medicine (BMC),
13 Department of Clinical Sciences, Lund University, Sweden; ⁶Department for BioMedical Research,
14 University of Bern, Switzerland; ⁷Department of Proteomics and Signal Transduction, Max Planck
15 Institute of Biochemistry, Martinsried, Germany; ⁸Institute of Biophysics, Chinese Academy of
16 Sciences, Beijing, China; ⁹Faculty of Science, University of Zürich, Switzerland.

17
18 †These two authors contributed equally to this work.

19 *To whom correspondence should be addressed.

20

21 **Abstract**

22 Protein complexes, macro-molecular assemblies of two or more proteins, play vital roles in numerous
23 cellular activities and collectively determine the cellular state. Despite the availability of a range of
24 methods for analysing protein complexes, systematic analysis of complexes under multiple conditions
25 has remained challenging. Approaches based on biochemical fractionation of intact, native complexes
26 and correlation of protein profiles have shown promise, for instance in the combination of size
27 exclusion chromatography (SEC) with accurate protein quantification by SWATH/DIA-MS. However,
28 most approaches for interpreting co-fractionation datasets to yield complex composition, abundance
29 and rearrangements between samples depend heavily on prior evidence. We introduce PCprophet, a
30 computational framework to identify novel protein complexes from SEC-SWATH-MS data and to
31 characterize their changes across different experimental conditions. We demonstrate accurate
32 prediction of protein complexes (AUC >0.99 and accuracy around 97%) via five-fold cross-validation
33 on SEC-SWATH-MS data, show improved performance over state-of-the-art approaches on multiple
34 annotated co-fractionation datasets, and describe a Bayesian approach to analyse altered protein-protein
35 interactions across conditions. PCprophet is a generic computational tool consisting of modules for
36 data pre-processing, hypothesis generation, machine-learning prediction, post-prediction processing,
37 and differential analysis. It can be applied to any co-fractionation MS dataset, independent of separation
38 or quantitative LC-MS workflow employed, and to support the detection and quantitative tracking of
39 novel protein complexes and their physiological dynamics.

40 **Main**

41 The analysis of proteins has progressed from studying specific proteins to the comparative analysis of
42 multiple proteomes, allowing for the detection of changes in the proteome landscape as a function of
43 the cellular state and the identification of connections between proteins based on their behaviours across
44 multiple samples. However, proteins largely function as complexes which are involved in performing
45 and regulating a majority of biological functions¹⁻⁴. Protein complexes are a part of extended functional
46 groups, such as pathways or protein interaction networks. Despite the availability of a range of methods
47 for the analysis of specific protein complexes, systematic analysis of the ensemble of protein complexes
48 in a sample has remained challenging⁵. While affinity-purification mass spectrometry (AP-MS)
49 provides valuable biological information on protein complexes, it lacks scalability and requires either
50 genetic manipulation of cells for introduction of a tag or the use of antibody-based reagents. On the
51 other hand, biochemical fractionation mass spectrometry allows for simultaneous quantification of
52 thousands of proteins and is emerging as a powerful technique for system-wide investigation of protein
53 complexes. Analytical techniques such as size exclusion chromatography (SEC) and ion exchange
54 chromatography (IEX) have been successfully applied in a variety of complex biological questions
55 such as apoptosis-dependent complex rewiring⁶, characterization of novel complexes in *Trypanosoma*
56 *Brucei*⁷ and *C. Elegans*⁸, identification of isoform-specific complexes⁹ and differential analysis of cell
57 cycle states¹⁷, i.e. the interphase and mitosis.

58 A key challenge in fractionation-based approaches is the confident assignment of protein
59 subunits to protein complexes based on their co-fractionation patterns and other relevant biological
60 information. A number of computational frameworks have been proposed for this purpose^{8, 10-12}.
61 Among these methods, CCprofiler identifies protein complexes from co-fractionation proteomic data
62 based on prior information from reference complex/interactome databases such as CORUM¹³,
63 STRING¹⁴ and BioPlex^{15, 16}. CCprofiler was not designed to predict novel protein complexes but to
64 determine a confidently detectable set of complexes including statistical estimation and control of the
65 false discovery rate¹⁰. PrInCE and EPIC leverage machine-learning techniques to predict novel protein
66 complexes but are limited conceptually to the inference of protein-protein interactions (PPI) from co-

67 fractionation proteomic data. Finally, dendrogram clustering has been described for novel complex
68 identification¹¹. In this case as well, control of false positives and false negatives is challenging since
69 an arbitrary threshold must be applied to cut the dendrogram¹¹.

70 In this study, we describe PCprophet, an open-access software for protein complex prediction
71 directly from co-fractionation-MS data using machine-learning techniques and differential analysis of
72 complex abundance and assembly state across conditions. PCprophet combines the benefits from
73 previous approaches such as error rate control using database-derived complexes present in CCprofiler
74 with the discovery of novel complexes inherent in other approaches. PCprophet offers the following
75 features: (i) PCprophet accepts input from a variety of co-fractionation mass-spectrometry (coFrac-MS)
76 techniques, including but not limited to size-exclusion chromatography (SEC-MS), strong cationic
77 exchange (SCX), and blue native page (BNP); (ii) PCprophet can be used with inputs derived from
78 widely employed mass spectrometry acquisition schemes such as data dependent acquisition (DDA),
79 data independent acquisition (DIA) and different quantitation strategies such isobaric labelling (SILAC,
80 TMT) or label-free; (iii) PCprophet was trained using co-eluting protein complex data, rather than co-
81 eluting PPIs, and can therefore directly predict novel protein complexes (i.e. complex-centric
82 prediction); (iv) PCprophet performs post-prediction processing via a statistical error model based on
83 Gene Ontology scores and other criteria to improve the reliability of the predicted protein complexes
84 and reduce false positives; (v) PCprophet performs differential analysis of predicted protein complexes
85 across conditions using our newly proposed Bayesian inference-based method. We applied PCprophet
86 to predict and analyse protein complexes in different cell cycle phases using our recently published
87 SEC-SWATH-MS dataset in the HeLa cell line¹⁷. Our results demonstrate that PCprophet predicts
88 novel protein complexes and recapitulates known changes in protein complexes across the cell cycle.

89

90 **Results**

91 **PCprophet accurately identifies novel protein complexes from co-fractionation MS data**

92 PCprophet enables accurate prediction of protein complexes directly from raw input (i.e. protein
93 matrices consisting of protein intensity vs. fraction number) of SEC-SWATH-MS and other co-

94 fractionation data. The framework of PCprophet (**Fig. 1**) includes six major modules: data pre-
95 processing, database query and *de novo* complex (i.e. hypothesis) generation, feature calculation and
96 prediction, error estimation and post-prediction processing, complex-centric differential analysis, and
97 report generation and data visualisation. During the data pre-processing step, Gaussian filtering,
98 missing value imputation, linear interpolation and data resizing are performed to ensure data quality
99 (See ‘**Methods**’ for more details). During the hypothesis generation step, a list of candidate protein
100 complexes for each condition based on the raw input protein matrices is provided separately via peak-
101 picking and distance-based clustering, for the machine-learning model to predict. During feature
102 calculation and prediction, each protein complex delivered by the hypothesis generation procedure is
103 represented using a numeric vector, including average intensity difference of proteins within each
104 fraction, local correlation of proteins at each window, shift of apex fraction of each protein and average
105 full width of a peak at half maximum. Meanwhile, the provided database (either PPI or complexes) is
106 mapped in the same feature space for later being used for FDR control. Then the Random Forest models
107 predict potential protein complexes with detailed predicted probabilities. During error estimation and
108 post-prediction processing, PCprophet filters the predictions based on Gene Ontology (GO) terms
109 assigned to components of predicted complexes. By calculating the pairwise GO term semantic
110 similarities of proteins assigned to a complex and comparing them to similarity scores in reference
111 databases of known protein complexes, PCprophet filters predicted complexes by a local false
112 discovery rate (FDR) based on GO term semantic similarity. In addition, PCprophet performs complex
113 combination and collapsing, since hypothesized complexes might be a subset of a bigger complex or a
114 mix of multiple complexes. During the complex-centric differential analysis, PCprophet analyses the
115 differences in prediction results between conditions, from protein level to complex level, using a
116 Bayesian inference method. As the final output, tabular and visual reports of the predicted protein
117 complexes and their changes across different conditions are generated by PCprophet. In summary,
118 PCprophet provides a ‘one-stop’ computational framework for the confident detection of protein
119 complexes including their dynamic changes across different biological states from a wide range of
120 coFrac-MS data.

121

122 **Benchmarking PCprophet complex prediction against state-of-the-art methods**

123 Concluding from the five-fold cross-validation (refer to ‘**Supplementary Results**’ for more details),
124 Random Forest (RF) has been chosen as the core classification algorithm of PCprophet. We then
125 assessed the performance of complex predictions using the optimized PCprophet framework against
126 two different, state-of-the-art computational approaches for the detection of protein complexes from
127 co-fractionation data, namely CCprofiler¹⁰ and EPIC⁸. Similar to PrInCE¹², EPIC supports the
128 prediction of binary protein-protein interactions and network inference of underlying complexes,
129 maintaining the potential to discover previously unknown complexes. Out of these two interaction-
130 centric approaches, we selected EPIC for our performance comparison as it has been shown to
131 outperform previous tools such as PrInCE. We benchmarked these tools based on a recently published
132 dataset, where HeLa CCL2 cells were synchronized in distinct cell-cycle stages (i.e. interphase and
133 mitosis) prior to analysis by SEC and DIA/SWATH-MS¹⁷. To avoid biases in assessing performance
134 due to the different inputs required by these tools (CCprofiler mainly takes the peptide-level
135 quantitative values as input, whereas PCprophet and EPIC take as input the protein-level quantitative
136 values), we performed sibling peptide correlation using CCprofiler and exported the resulting protein
137 matrices, thereby providing the same input for all benchmarked tools (refer to ‘**Methods**’ for more
138 details). To minimize comparison bias due to parameter optimization, we ran CCprofiler with the
139 parameters used in its original publication¹⁰. EPIC, on the other hand, offers the possibility of choosing
140 between an SVM classifier or an RF classifier for PPI prediction. We used default parameters with both
141 classifiers, generating two sets of predictions (EPIC_SVM and EPIC_RF). We generated protein
142 complex hypotheses for CCprofiler using the CORUM core complexes dataset, and also trained EPIC
143 using CORUM. PCprophet requires a protein complex or PPI database as input to perform FDR control
144 and CORUM was used for this purpose as well.

145 We initially evaluated the performance of each method using the numbers of known CORUM
146 complexes recovered across all replicates and conditions. Both the absolute number of identified
147 complexes as well as the overall recall are vastly different for each tool. The complex-centric tools (i.e.

148 PCprophet and CCprofiler) identified 900 and 798 known complexes respectively; while EPIC_RF
149 recovered only 71 known complexes and EPIC_SVM recovered none (**Fig. 2a**). The overlap in known
150 complexes between PCprophet and CCprofiler was 69.7% (i.e. 556 out of 798), while 49.2% (i.e. 35
151 out of 71) overlap was achieved between PCprophet and EPIC_RF (**Fig. 2b**). The identified complexes
152 correspond to a recall rate of 37% for PCprophet, 33% for CCprofiler and 3% for EPIC_RF. PCprophet
153 recalls a much higher fraction of CORUM complexes compared to those recalled by EPIC analysis also
154 in a DDA-based dataset with a isotope dilution strategy for quantification¹⁸ (DDA-SILAC,
155 **Supplementary Fig. S1**). We then compared the average number of subunits per complex to evaluate
156 the similarity of known complexes from CORUM to the predicted ones from EPIC and PCprophet (**Fig.**
157 **2c**). The distribution of subunits per complex predicted by PCprophet and CCprofiler is closer to that
158 of CORUM complexes (average 4.1 subunits), with an average subunit size of 3.5 for PCprophet and
159 6.9 for CCprofiler. The average subunit size per complex predicted by EPIC, however, was 19.9
160 ($p < 10E-14$) for the RF classifier and 71.7 ($p < 10E-14$) for the SVM classifier, respectively. The results
161 from EPIC thus suggest a larger size of cellular assemblies compared to the sizes of manually curated
162 complexes in the CORUM database, with more similar sizes reported by both CCprofiler and
163 PCprophet.

164 We then evaluated the performance of PCprophet and EPIC in recalling protein-protein
165 interactions (PPIs). In this comparison, we did not consider CCprofiler as it cannot derive novel
166 complexes without prior information, which limits its applicability for discovery of novel protein-
167 protein interactions. We generated a PPI network from complexes predicted by PCprophet, EPIC_RF,
168 and EPIC_SVM, and compared them to ground truth networks from CORUM complexes and from PPI
169 databases such as STRING and BioPlex. This comparison allows to calculate the percentage of reported
170 PPIs for each tool, in the form of PPI precision. PCprophet achieved a PPI precision of 0.65 when
171 compared with STRING and 0.095 in comparison to BioPlex database (**Fig. 2d**). The precision of
172 EPIC_RF was 0.12 with STRING and 0.004 when compared with PPIs in BioPlex. EPIC_SVM
173 prediction corresponded to a precision of 0.11 and 0.002 with STRING and BioPlex respectively. We
174 calculated for each network the degree distribution and the frequency of nodes with a particular degree

175 **(Fig. 2e)**. To evaluate the similarity between ground-truth networks and prediction, the Area Under
176 the Curve (AUC) values were calculated for all the tools (**Supplementary Table S1**) and databases.
177 Regardless of the classifier used, EPIC-derived PPI networks tend to have higher degree (**Fig. 2e**)
178 compared to those from complexes in CORUM. This resulted in an AUC of 0.18 for EPIC_RF, 0.39
179 for EPIC_SVM, 0.13 for PCprophet and 0.11 for CORUM, respectively. In this context, an AUC value
180 closer to the one of reported complexes (CORUM) means a closer resemblance in network topology to
181 a ground truth network. Finally, we merged all PPI databases (STRING, BioPlex and BioGrid) to
182 generate a combined network including all deposited interactions and assessed the average distance
183 between every pair of proteins within a predicted or known complex. The average shortest path for
184 EPIC_RF was 2.3 and >3 for EPIC_SVM while PCprophet-predicted complexes had an average path
185 of 1.1 edges as shown in **Fig. 2f**, suggesting a greater recovery of closely connected proteins by
186 PCprophet when compared to the average shortest path in CORUM (1 edge). We also observed a
187 similar trend on an independent dataset from PrInCE¹² (**Supplementary Fig. S1**). To summarize,
188 PCprophet allows for robust identification of complexes, as shown based on high recall of known
189 complexes and high quality of newly predicted complexes, demonstrated based on the high validation
190 rates of the underlying PPIs by large-scale databases. PCprophet outperforms available tools in the
191 recovery of known protein complexes and PPIs while additionally providing the opportunity to detect,
192 investigate and track assemblies that remained inaccessible to computational approaches limited by
193 their dependence on prior knowledge¹⁰ or the sensitivity of interaction-centric scoring^{8, 12, 19}. In order
194 to control spurious co-elution and false positive assignments, we integrated an error model based on
195 interactor gene ontology similarity which effectively ensures highest quality of the reported results
196 (**Supplementary Fig. S2**).

197

198 **Predicting PPIs and protein complexes across the mammalian cell cycle via PCprophet**

199 We applied PCprophet to a second, newly published dataset¹⁷ in which HeLa cells were blocked at
200 mitosis and interphase stages of the cell cycle. Proteins were then extracted under native conditions,

201 SEC separated into 65 fractions and analysed using SWATH-MS. Based on these data, we generated a
202 large PPI map based on all PCprophet predictions (**Fig. 3a**).

203 PCprophet predicted 858 protein complexes not recorded in the CORUM derived network,
204 which contain 11527 unique PPIs consistently present across all biological replicates for one condition
205 (**Fig. 3b**), suggesting good reproducibility across different fractionation experiments. Of these
206 predicted PPIs, 54.16% are consistently supported by evidence across several PPI databases (STRING,
207 BioPlex and BioGrid²⁰), 14.67% have PPI evidence from a single database while 31.14 % of the PPIs
208 are completely novel (**Fig. 3c**), consistent with the 30% FDR cut-off used for the search (refer to
209 ‘**Methods**’ section for more details). FDR in PCprophet is calculated by comparing hits from the
210 provided database, in this instance CORUM, against positively predicted complexes, thereby 30% of
211 the PPIs detected cannot be derived from CORUM. We speculated that this set of PPIs without
212 database evidence would be localized in cellular niches with poorly characterized complexes, such as
213 membrane bounded organelles⁸. Consistent with this hypothesis, among the top 10 most enriched GO
214 Cellular Compartments (CC) terms for the novel PPIs, we observed localization enrichment in
215 mitochondrion, ficolin, cytoskeleton and cytoplasmic-associated lumen (**Fig. 3d**) all with an adjusted
216 *p*-value of less than 1%.

217 We identified several cases where a novel subunit is assigned to a known complex by PCprophet.
218 For instance, PCprophet identified a novel protein complex containing the ubiquitin receptor ADRM1
219 and 26S proteasome (**Fig. 3e**). This association has not been reported in the CORUM database for *homo*
220 *sapiens* but it has been identified in mammalian cells²¹ and is consistent with the crystal structure of *S.*
221 *cerevisiae* (PDB ID: 6J2C and 6J2Q)²². ADRM1 is reported to be a component of the 19S proteasomal
222 subunit in yeast²²; accordingly we observed about 15% of the ADRM1 signal to be associated with 19S
223 (**Fig. 3e**) while the majority was associated with the 26S proteasome, suggesting that ADRM1 is
224 preassembled in the 19S rather than being later recruited to a fully assembled 26S. We further identified
225 an interaction between the NEDD8 activating complex NAE1-UBA3 and ASB6 (**Fig. 3f**). The NAE1-
226 UBA3 complex is required for cell cycle progression by transferring activated NEDD8 to UBE2M and
227 subsequent proteasomal degradation²³. ASB6, which belongs to the Ankyrin repeat and SOCS box

228 (ASB) protein family, has been shown to interact with CUL5 and RBX2 to form a non-canonical E3
229 ubiquitin ligase complex²⁴. We observed almost perfect co-elution between the NAE1-UBA3 complex
230 and ASB6, consistent with reports of ASB6 and UBA3 co-purification in other species^{25, 26}, but not
231 with reported ASB6 binders such as CUL5²⁴ (**Supplementary Fig. S3**) or reported NAE1-UBA3
232 binders like UBE2M¹⁶ or TP53BP2¹⁶. Taken together, the recall of protein-protein interactions absent
233 from the training set as well as reported complexes, suggests that PCprophet can predict protein
234 complexes in cellular models that are poorly characterized with respect to protein complexes and PPIs.

235

236 **Differential analysis of mitosis-associated protein complexes**

237 We have identified 900 previously reported and 532 novel complexes in HeLa cell lysates derived from
238 interphase and mitotic cells, with a similar number of complexes in each cellular state (**Fig. 4a**). Due
239 to the continuous nature of coFrac-MS data it is possible to identify several types of profile differences
240 at both protein and complex level. First, difference in assembly state causes a shift on the molecular
241 mass scale, while changes at the abundance level results in an increase peak area for a particular protein.
242 By similarity, complex compositional changes can be inferred by the difference in peak position of the
243 subunits or addition of novel proteins, while stoichiometric changes are dependent on ratios between
244 different proteins. This is a non-trivial issue as metrics based on profile correlation will fail in capturing
245 abundance difference, while methods based on peak position will not detect variation in peak area. To
246 overcome this issue, we developed a Bayesian approach to identify altered protein profiles in the
247 different conditions tested and defined a likelihood for each interaction, which we then combine into a
248 complex-specific likelihood (see '**Methods**' for more details). This approach has several advantages
249 over previous methods such as fold change¹⁷ as it does not require a pre-selected threshold and
250 penalizes proteins with high variability. Overall, we detected 1518 proteins (238 complexes) with a
251 probability greater than 0.5 of being differentially regulated across the cell cycle (**Fig. 4b**). On this set
252 of proteins, we performed an enrichment analysis using GO ontology to evaluate if terms associated
253 with cell cycle and mitosis (**Fig. 4c**) were enriched. Indeed, terms such as M phase, mitosis and nuclear
254 division are enriched with $p < 0.001\%$. Surprisingly, we identified the Prmt5-Wdr77 complex as altered

255 between interphase and mitosis (**Fig 4d**). This complex is composed by a hetero tetramer formed by
256 Prmt5-Wdr77 dimers in a 1:1 ratio²⁷ (PDB ID: 4GQB). While this putative 1:1 stoichiometric ratio is
257 reflected in the protein MS intensities during interphase, upon mitosis it is significantly shifted towards
258 a 1.75:1 ratio, indicating a gain of Prmt5 copies in the assembly relative to the composition in interphase
259 (**Fig 4e**). Interestingly, Prmt5²⁸ and Wdr77²⁹ have been independently linked to cell cycle regulation
260 and complex stoichiometry is necessary for correct target methylation by Prmt5²⁷. Thus, our data
261 suggests a potential role for the Prmt5-Wdr77 complex in cell cycle regulation. Furthermore, key events
262 such as activation of the master mitotic kinase complex CDK1/CCNB1 (**Fig. 4f**), increase in cohesin
263 complex (**Fig. 4g**) and rewiring of the anaphase promoting complex/cyclosome (**Fig. 4h**) were
264 successfully captured by our analysis strategy.

265 To conclude, our analysis demonstrates that (i) its ability to recall known complex remodelling
266 events in cell cycle progression, and (ii) its sensitivity to discriminate between different scenarios such
267 as increase in abundance (**Fig. 4fg**) and difference in peak shapes (**Fig. 4h**). Altogether, our analysis
268 recapitulates previous knowledge about cell cycle and cell cycle-related events, selectively recalling
269 complexes involved in cell cycle progression and mitosis.

270

271 **Discussion**

272 Protein complexes play fundamentally important roles in mediating and regulating biological functions.
273 Recent advances in proteomic technologies based on co-fractionation and mass spectrometric
274 correlation profiling of protein elution patterns have opened up a promising avenue to characterize
275 protein complexes at breadth and temporal resolution. State-of-the-art workflows such as SEC-
276 SWATH-MS techniques and complex-centric data analysis have advanced the selectivity and
277 throughput of chromatographic protein complex detection but remain limited to the detection of
278 previously observed protein complexes. Methods to predict novel protein complexes from co-
279 fractionation data are based on identification of PPIs and inference of complexes from the resulting
280 weighted network. Such probabilistic methods for network partitioning rely heavily on network
281 topology, which makes it challenging to partition detected PPIs into complexes, due to the high

282 dimensionality of the data. In light of this, we introduced the PCprophet framework which combines
283 complex-level scoring with powerful machine learning technology to classify and confidently predict
284 novel protein complexes from protein coFrac-MS data. In addition, PCprophet facilitates the
285 differential tracking and comparison of these complexes across two or more experimental conditions
286 that become increasingly accessible via high throughput implementations of coFrac-MS. We have
287 demonstrated outstanding prediction performance of PCprophet on manually annotated datasets and
288 have shown that the method significantly outperforms state-of-the-art complex prediction and
289 identification tools. We have developed a Bayesian inference-based method to analyse differences in
290 protein complex abundance and composition across conditions. Our analysis on proteomic profiles
291 across the cell cycle of HeLa cells demonstrated that PCprophet can capture expected changes in
292 protein complexes between interphase and mitotic cells.

293 PCprophet is available in command-line version under MIT licence
294 (<https://github.com/fossatiA/PCprophet>) and is easily applied to any coFrac-MS dataset. The data pre-
295 processing module readily accepts different types of quantitative protein level tables. PCprophet could
296 also be applied in clinical proteomics and personalized medicine areas, to assist the discovery and
297 analysis of novel protein complexes and to identify complexes that are altered across groups of samples.
298 We anticipate that the PCprophet package will serve as a reliable and accurate tool for novel protein
299 complex prediction and analysis from co-fractionation MS data because it extends the scope of
300 comparative proteomics from the level of differentially abundant proteins to the level of differentially
301 abundant and perturbed complexes between samples, thus bringing proteomic analysis closer to
302 biological function.

303

304 **Methods**

305 Methods, including statements of data availability and any source code and references, are available in
306 the online version of the paper.

307

308

309 **Acknowledgements**

310 This work was supported in part by the Swiss National Science Foundation (grant No. 3100A0-688
311 107679 to R.A.) and the European Research Council (ERC-2014AdG 670821 to R.A.). C.L. is
312 currently supported by a National Health and Medicine Research Council (NHMRC) of Australia CJ
313 Martin Early Career Research Fellowship (1143366). M.G. and F.F. acknowledge the support by the
314 Innovative Medicines Initiative project ULTRA-DD (FP07/2007-2013, grant No. 115766). I.B.
315 acknowledges funding support from the Swiss National Science Foundation (31003A_166435). AWP
316 is supported by a NHMRC Principal Research Fellowship (1137739). The authors would like to thank
317 Dr. Natalie de Souza from ETH Zurich and Associate Professor Jiangning Song from Monash
318 University for their critical comments on this study.

319

320 **Author contributions**

321 R.A., A.F., C.L. and M.G. conceived and designed the project. A.F. and C.L. designed, developed and
322 implemented PCprophet, and conducted data analysis, machine-learning prediction and benchmarking
323 experiments with other existing methods. P.S. designed and implemented the differential analysis
324 module for protein complexes. M. Heusel, F.F. and F.U. contributed to data annotation and provided
325 critical feedback and comments on the biological aspects. M. Hallal contributed to EPIC performance
326 comparison and reproducibility test for PCprophet performance. I.B. assisted with benchmarking
327 experiments with CCprofiler and provided useful insights. C.T.K., P.X. and A.W.P. provided critical
328 and insightful comments during the development of PCprophet. C.L., A.F., M.G. and R.A. drafted the
329 manuscript, which has been revised and approved by all the other authors.

330

331 **Competing financial interests**

332 The authors declare no competing financial interest.

333 **Methods**

334 **Training dataset curation and annotation**

335 In total, three co-fractionation replicates using SWATH and DDA-SILAC based datasets were used to
336 train and evaluate PCprophet, including the SEC-SWATH-MS dataset from HEK293 cell line¹⁰ for
337 training PCprophet, the mitotic proteomic data from HeLa CCL2 cells¹⁷ and the DDA-SILAC dataset
338 extracted from the study by Kristensen and used as testing dataset for PrInCE¹², for independently
339 testing PCprophet. Note that we did not use the *C. elegans* protein complex dataset from the EPIC⁸
340 package to test PCprophet for the following reasons: (i) their datasets used spectral count; however our
341 previous study showed lower performance of spectral counts compared to XIC based quantitation (MS1
342 or MS2) for complex analysis¹⁰; (ii) the features and pre-processing employed in PCprophet are
343 inherently continuous in nature such as correlation and FWHM; and (iii) EPIC was developed and
344 evaluated using the same dataset. It is therefore challenging to conduct an unbiased and fair estimation
345 and comparison of the performance for all the other approaches. A structuralized description of these
346 datasets is available in **Supplementary Table S2**.

347 To train accurate machine-learning models, we manually annotated the protein complexes from
348 the SEC-SWATH-MS dataset from HEK293 cells¹⁰. Briefly, samples were acquired in SWATH mode
349 using a sample-specific library generated from high-pH fractionated samples. Following conversion to
350 mzXML via msConvert, OpenSWATH search was performed with the parameter previously
351 described¹⁰. As a result, the final feature alignment outputs from TRIC³⁰ (using top2 protein
352 quantification) was used for the training PCprophet. The reference core (non-redundant) complexes
353 were downloaded from CORUM v3.0¹³. Protein accession numbers were converted into gene names
354 and sequentially mapped to CORUM. We removed complexes from our dataset where the number of
355 subunits present in the dataset was less than 50% of known components to retain only complexes with
356 high coverage, consistent with the annotation strategy for the complex analysis software CCprofiler¹⁰.
357 To train a supervised classification model, it is crucial to reliably annotate the samples of positive (i.e.
358 complexes with good co-elution profiles) and negative (i.e. complexes with poor co-elution profiles)
359 classes. For the training dataset (i.e. HEK293), a protein complex was annotated positive if it satisfied

360 the following criteria. First, more than 75% of the known subunits coeluted in the same fraction with
361 baseline resolution; second, the main peak was required to have a minimal FWHM (full width at half
362 maximum) of 4 fraction; third, minimal normalized height of 20% to the maximum signal for every
363 protein and needs to be at least 10% above background, and the complex has been annotated in the
364 CORUM database. On the other hand, if the complex was annotated in CORUM but did not pass all
365 the other criteria it was annotated as negative. To objectively annotate the protein complexes in the
366 training dataset, three annotators were involved in this procedure and only positive and negative protein
367 complexes nominated and agreed by all the annotators were used. Notice that we did not randomly
368 select proteins to form negative but ‘fake’ protein complexes, as there would be a huge number of
369 different combinations and possibilities and might result in random selection of undiscovered protein-
370 protein interaction. In addition, we changed the original number of fractions of the training dataset from
371 81 to 72, to standardize number of fractions across different experiments. The final resulting training
372 dataset contained 242 positive protein complexes and 738 negative complexes.

373

374 **Dataset pre-processing**

375 Prior to the generation of potential protein complexes based on the protein raw matrices across various
376 conditions, four data pre-processing steps are performed, including Gaussian filtering, missing value
377 imputation, linear interpolation, and data rescaling. One-dimension Gaussian filtering,

$$G(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}} \quad (1)$$

378
379 was employed to smooth the data by removing noise and approximating peaks as Gaussian curves,
380 where x denotes the intensity of the current fraction and σ was set to 1. To remove the missing values
381 in the raw data, an imputation strategy by calculating the average of the two neighbours of missing
382 value, was implemented. The number of fractions N in the co-fractionation experiments always varies
383 due to various experimental setups and inherent variability. When constructing the machine-learning
384 models, the number of features is dependent on the number of fractions (see ‘Feature engineering and
385 construction of machine-learning models’ for feature generation). Based on the number of fractions

386 (i.e. 72) in our training dataset (see ‘Data curation and annotation’ for details), it was therefore
387 necessary to rescale the number of fractions of user-provided datasets to 72. In PCprophet, resampling
388 and one-dimensional linear interpolation were applied for this purpose, thereby rescaling the number
389 of fractions to 72, consistent with the training dataset. Lastly, we added an additional step to standardize
390 every protein profile from their original intensity in the range [0, 1] to make it independent from the
391 quantitation strategy used.

392

393 **Hypothesis generation**

394 In this study, hypothesis generation refers to the construction of putative protein complexes.
395 Theoretically, there could be a huge number of possible combinations of proteins to form different
396 protein complexes. Rather, we proposed a hypothesis generation module to construct potential protein
397 complexes by aligning peaks of different proteins and cutting the dendrogram-like tree structure,
398 similar to the procedure discussed in a previous study¹¹. To do so, hypothesis generation firstly
399 performs peak-picking to identify all the apexes of intensity and their associated fractions of all proteins
400 in the input data, with the help of the Python package ‘SciPy’ (<https://www.scipy.org/>). Then linkage
401 hierarchical clustering using the ‘Ward’ distance measure was performed based on the apexes collected
402 during the peak-picking stage. A dendrogram-like tree structure was then generated based on the results
403 of the linkage hierarchical clustering. The hypotheses (i.e. putative protein complexes) were then
404 generated by cutting the dendrogram from bottom to top at each level. A conceptual illustration of the
405 hypothesis generation is presented in **Supplementary Fig. S4**. In practice, these steps are performed
406 on a linkage matrix instead of a dendrogram structure to further reduce the computational burden.

407

408 **Feature engineering and construction of machine-learning models.**

409 To represent a protein complex, we designed a variety of features based on the protein co-elution profile
410 and the number of fractions N . These features are mainly categorised in four groups: (1) average
411 intensity difference of proteins within each fraction (**Supplementary Fig. S5a**), (2) local correlation

412 of proteins at each window (**Supplementary Fig. S5b**), (3) shift of apex fraction of each protein
413 (**Supplementary Fig. S5c**), and (4) average full width half maximum (**Supplementary Fig. S5d**).

414 For the intensity difference of proteins and the correlation of proteins at each fraction, we set a
415 sliding window (6 fractions wide) with 1 fraction step wise increase. The intensity difference of proteins,
416 $D_{intensity_i}$ reflects the average difference of intensity values of protein a , b , and c , at each fraction,
417 calculated by:

$$D_{intensity_i} = \mu(|a_i - b_i|, |a_i - c_i|, |b_i - c_i|), i = (1, 2, \dots, N) \quad (2)$$

418 where i denotes the number of fractions. As a result, the dimension of this feature type is N . Similarly,
419 at each window, pairwise correlation of intensity was also calculated, using:

$$Corr_i = \mu(Corr(a_{wi}, b_{wi}), Corr(a_{wi}, c_{wi}), Corr(b_{wi}, c_{wi})), i = (1, 2, \dots, N) \quad (3)$$

421 Where a_{wi} denotes the local value of a in a window w centred at fraction i . For the other two types of
422 features, including fraction difference of apex peaks and full width half maximum a two step-procedure
423 is employed. First, all peaks for a protein complex hypothesis are selected, and then a modified version
424 of the Dijkstra's algorithm is applied to select the peaks for every protein with the minimum distance.
425 By selecting the closer peaks, we are able to positively predict proteins with multiple peaks in separate
426 assemblies, as we avoid the use of heuristic to select the complex-specific peak. While the apex
427 difference is a feature used also in the PrInCE software¹², substantial differences are present as in this
428 tool, the fraction with the maximum value for every protein is counted as apex, thereby using always
429 the same peak for a protein in multiple assemblies.

431 For the average apex difference, we used following formula for the calculation, respectively:

$$F_{a,b,c} = \mu(|X_{pa} - X_{pb}|, |X_{pa} - X_{pc}|, |X_{pb} - X_{pc}|), \quad (4)$$

432 where X_{pa} , X_{pb} , X_{pc} represent the apex fraction of protein a , b and c , respectively. The average full
433 width half maximum is calculated using:

$$FWHM_{a,b,c} = \mu(|X_{a_right} - X_{a_left}|, |X_{b_right} - X_{b_left}|, |X_{c_right} - X_{c_left}|), \quad (5)$$

435

436 while X_{a_right} - X_{a_left} , X_{b_right} , X_{b_left} , X_{c_right} , X_{c_left} demonstrate the width when achieving half
437 intensity area of the co-elution curve of protein a , b and c , respectively. In total, we generated $2N + 2$
438 features (i.e., 146 when $N=72$).

439 Five well-established machine-learning models were selected to test the prediction performance,
440 including Decision Tree (J48)³¹, Random Forest³² (RF), Naïve Bayes³³ (NB), Support Vector
441 Machines³⁴ (SVM) and Logistic Regression³⁵ (LR). For SVM, we selected two major kernels, including
442 polynomial and RBF³⁶ (Radial Basis Function) kernels, due to the consideration of the balance of
443 computational complexity and running time. These two models were then termed as SVM_POLY and
444 SVM_RBF, respectively. Note that the above machine-learning algorithms were implemented using
445 the scikit-learn package³⁷ and cross-tested in the WEKA³⁸ platform. Different implementations of such
446 algorithms in other platforms may cause difference in term of prediction performance. To objectively
447 portrait the prediction performance and avoid overfitting, five-fold cross-validation strategy using the
448 training dataset was performed, together with five widely acknowledged performance measures,
449 including accuracy (ACC), area under the curve (AUC), Matthew's correlation coefficient (MCC),
450 sensitivity and specificity:

$$Acc = \frac{TP + TN}{TP + TN + FP + FN} \quad (6)$$

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (7)$$

$$Sensitivity = \frac{TP}{TP + FN} \quad (8)$$

$$Specificity/TPR = \frac{TN}{TN + FP} \quad (9)$$

455 where TP , TN , FP , FN are true positives, true negatives, false positives and false negatives, respectively.

456

457

458

459 **Benchmarking with state-of-the-art approaches for co-fractionation MS based protein complex** 460 **prediction**

461 We compared the prediction performance of PCprophet with currently existing computational
462 approaches for protein complex characterization based on co-fractionation MS data, including
463 CCprofiler¹⁰, EPIC⁸, and PrInCE¹². CCprofiler is a statistical approach for the identification of protein
464 complexes by referencing databases as prior information. In contrast, EPIC and PrInCE were designed
465 to infer protein complexes based on PPI prediction from co-elution profiles. When running EPIC, two
466 provided models, including SVM and RF (i.e. -M RF and -M SVM) were both tested with other
467 parameters by default (i.e. -t 9606 for *H. sapiens*; -s 11101001; -f STRING). Two datasets, as shown
468 in **Supplementary Table S2** were used to benchmark with CCprofiler and EPIC, including the mitotic
469 proteomic data from HeLa cells¹⁷, and the soluble protein complex dataset from the study of Stacey *et*
470 *al.*¹²

471 *Benchmarking using HeLa mitotic proteomic data.* The TRIC feature aligned file of the dataset
472 was imported into CCprofiler and following sibling peptide correlation (the ‘filterBySibPepCorr’
473 function). Protein quantification was done using the top 2 proteotypic peptides per protein. The
474 resulting protein tables were exported and used as input for EPIC and PCprophet. PCprophet was run
475 with default parameters, with the FDR fixed at 30% and controlled using the CORUM database.
476 CCprofiler complex-centric analysis was done as previously described¹, using `smoothing_length = 9`,
477 `corr_cutoff = 0.95`, `window_size = 8`, `rt_height = 3` and a 2x molecular weight cutoff¹. For the
478 calculation of recall against CORUM, PCprophet output (i.e. the ‘ComplexReport.txt’ file) was filtered
479 to only ‘Reported’ complexes which were predicted as ‘Positive’; while EPIC derived complexes were
480 matched to CORUM by defining a positive predicted complex in which 50% or more subunits are
481 reported in a single CORUM core complex. For CCprofiler, the positive complexes were defined when
482 the q-value is smaller than 5%, as previously described¹. The number of complexes was considered
483 across replicates as the number of unique positive CORUM Complex ID. Two measures were applied
484 when assessing the prediction performance, including True Positive Rate (TPR; specificity) and
485 Positive Predicted Value (PPV), which is defined as follows:

$$PPV = \frac{TP}{(TP + FP)} \quad (10)$$

486
487 For evaluation of average network degree, protein complexes from the prediction outputs of EPIC and
488 PCprophet, and CORUM (Human only) were collapsed to a PPI network. The reference networks from
489 STRING (Human) and BioPlex were downloaded and used directly. Degree calculation for every
490 protein in the networks was done using the NetworkX package v2.1 (<https://networkx.github.io>) and
491 ranked. A log-log plot was generated and the AUC for the resulting curve was calculated using the
492 integrate module from SciPy (<https://www.scipy.org>). For evaluation of node centrality, the resulting
493 complexes from PCprophet, EPIC_RF and EPIC_SVM were projected into a subgraph generated by
494 filtering STRING to only nodes present in the original protein matrixes thereby representing all the
495 reported protein-protein interaction available in our data. For every complex in the three tools
496 (PCprophet, EPIC_RF and EPIC_SVM) the average complex closeness (ACC) was defined as the
497 mean shortest path between all members. The resulting vector represents the tendency of the different
498 algorithms to recapitulate first, second or outer shells level of interactors. The same was done also for
499 CORUM complexes within the same subgraph. Average complex size was defined as the mean number
500 of subunits for the same complex across replicate for PCprophet and the number of subunits for every
501 complex in the EPIC output.

502 *Benchmarking using the DDA-SILAC dataset*¹⁸. Condition1.tsv and condition2.tsv were
503 downloaded from <https://github.com/fosterlab/PrInCE-Matlab> and separated into the different
504 replicates. PCprophet and EPIC were run as described above. TPR and PPV were calculated as
505 described above for CORUM. For BioPlex and STRING both networks were filtered to remove
506 proteins not identified. PCprophet, EPIC_RF and EPIC_SVM derived complexes were collapsed to
507 generate a PPIs network and then recall was calculated using a PPI-centric approach by assessing which
508 fraction of PPIs was present over the entire reference (TPR) and which fractions of PPIs was correct
509 across all of the predicted one (PPV). Centrality assessment via KS test and AUC calculation was done
510 as described above.

511

512

513 **Post-prediction processing**

514 During this stage, GO (Gene Ontology) term score filtering and complex combination and collapsing
515 are performed in order to ensure the reliability of predicted complexes. Despite the incompleteness of
516 GO term annotation of CORUM database, we compared the distributions of GO terms of predicted
517 protein complexes and documented protein complexes in the CORUM database¹³ (**Supplementary Fig.**
518 **S6a**). We first collected GO terms of each protein in a predicted complex based on the annotation of
519 AmiGO2 database (i.e. the Gene Ontology resource)³⁹⁻⁴¹. Then, for every possible protein pair we
520 calculated pairwise GO term semantic similarity using the strategy published by Wang *et al*⁴². For
521 instance, given a protein complex PC with three subunits A , B and C , all the GO terms, including
522 molecular function (MF), biological process (BP), and cellular component (CC) are collected. For all
523 the three possible protein pairs, including A - B , A - C and B - C , within each category, the semantic
524 similarity scores of all the pairwise GO terms are calculated and the average score is reported as the
525 overall score for the current GO category. The final overall GO score of protein complex PC in this
526 case is then defined as follows:

$$GO(PC) = \mu(MF(A, B), MF(A, C), MF(B, C)) + \mu(BP(A, B), BP(A, C), BP(B, C)) \\ + \mu(CC(A, B), CC(A, C), CC(B, C)). \quad (11)$$

527

528 The GO term scores of core protein complexes from the CORUM database are calculated using the
529 same strategy. Given the two distributions of the GO term scores of both protein complex hypothesis
530 and the complexes harboured from the CORUM database, we then estimated false discovery rate for
531 the positively predicted hypothesis by calculating global FDR for every GO scores of positive CORUM
532 complexes with the following formula:

$$FDR_{(s|GO_h, GO_c)} = \frac{\int_s^\infty GO_h(x) dx}{\int_t^\infty GO_c(x) dx}, \quad (12)$$

533

534 where GO_h and GO_c are defined as distributions of the Wang similarity score for hypothesis (h) and
535 CORUM-derived complexes (c). This allows us to obtain the specific GO term score that satisfies the

536 target FDR for filtering the predicted protein complexes without having to use a fixed threshold, thereby
537 allowing for more or less conservative searches.

538 Given the possibility that the positive complex hypothesis might be a subset of a bigger complex
539 or might contain multiple smaller complexes, PCprophet allows users to select a complex collapsing
540 mode to further process the predicted complexes (**Supplementary Fig. S6b**), including ‘GO’ (based
541 on GO terms), ‘CAL’ (based on the provided calibration curve), ‘SUPER’ (to find the biggest protein
542 complex), and ‘NONE’ (to ignore this process). Specifically, following prediction and FDR control,
543 overlap is calculated and complexes for which the overlap is more than 0.75 are merged on the different
544 criterion. Given a set of complexes defined as follows in **Supplementary Table S3**, for example,
545 collapsing using ‘GO’ will result in the complex which has the greatest GO score (i.e. **PC2**). The
546 ‘SUPER’ option will select the of the complex with the highest number of subunits (i.e. **PC3**) while
547 choosing ‘CAL’ will calculate the difference between apparent MW from the SEC and extrapolated
548 molecular weight from the calibration curve. The complex with the smaller difference (i.e. **PC3**) will
549 therefore be selected. The ‘CAL’ mode is selectable only if the calibration and a molecular weight
550 table from the UniProt⁴³ database or similar format is provided. ‘NONE’ option will skip the collapsing
551 procedure.

552

553 **Protein complex differential analysis across different conditions**

554 *Bayesian inference of differential regulation of protein abundance.* Inferring differentially regulated
555 proteins assumes that protein abundance measurements were obtained for a number of samples which
556 differ in a biological phenotype of interest. This situation allows representing the protein abundance
557 measurements for one protein as a matrix X where rows correspond to samples and columns correspond
558 to retention time. Phenotype information t is assumed to be discrete with cardinality $\#t$ and order
559 matched such that the phenotype information for sample n , $t_n = t[n]$ corresponds to the protein
560 abundance row vector $x_n = X[n]$. If we assume that the correct model is among the investigated
561 candidate models, we have a problem termed ‘m-closed model selection’⁴⁴. In this situation the
562 Bayesian approach to inferring whether a phenotype change corresponds to differential regulation in

563 protein abundance suggests to use marginal likelihoods to derive the corresponding Bayes factors⁴⁴. To
 564 obtain a solution which may be calculated analytically, we use the model illustrated in **Supplementary**
 565 **Fig. S7a**. The model represents protein abundance measurements by Normal-Whishart distributions.
 566 Differential regulation is implicitly represented (variable not shown in the graph) via a protein specific
 567 indicator variable I_p . Differential regulation is coded by $I_p = 1$ and results in modelling the protein
 568 abundances x_n conditionally on phenotype states t_n by phenotype specific Gaussian distributions. To
 569 obtain a measure of differential protein regulation we compare the $I_p = 1$ model with a simpler
 570 explanation which we denote as $I_p = 0$. The simple model corresponds to a non-differential regulation
 571 assumption and uses one common Gaussian distribution to model x_n irrespective of the phenotype state
 572 t_n . Irrespective whether we have $\#t$ multivariate Gaussians in case of $I_p = 1$ or one shared Gaussian
 573 in case of $I_p = 0$, the joint distribution of data and model parameters $P(\mathbf{X}, \mu, \Lambda | \mathbf{t}, \gamma, m, g, h)$ is
 574 represented by the directed acyclic graph (DAG) in **Supplementary Fig. S7a**. In case of $I_p = 1$ we
 575 have

$$p(\mathbf{X}, \mu, \Lambda | \mathbf{t}, \gamma, m, g, h, I_p = 1) = p(\Lambda | g, h) \prod_{\tau=1}^{\#t} p(\mu_{\tau} | m, \gamma, \Lambda) \prod_n p(x_n | \mu_{t_n}, \Lambda), \quad (13)$$

576 whereas in case of $I_p = 0$ we have the simpler relation

$$p(\mathbf{X}, \mu, \Lambda | \mathbf{t}, \gamma, m, g, h, I_p = 0) = p(\Lambda | g, h) p(\mu | m, \gamma, \Lambda) \prod_n p(x_n | \mu, \Lambda). \quad (14)$$

577 The next step to obtain a measure of differential regulation is to calculate the marginal likelihood for
 578 both models in Equation (13) and Equation (14). For Equation (13) we get

$$p(\mathbf{X} | \mathbf{t}, \gamma, m, g, h, I_p = 1) = \int_{\Lambda} \int_{\mu_{\tau} \forall \tau} \left(p(\Lambda | g, h) \prod_{\tau=1}^{\#t} p(\mu_{\tau} | m, \gamma, \Lambda) \prod_n p(x_n | \mu_{t_n}, \Lambda) d\Lambda \prod_{\tau=1}^{\#t} d\mu_{\tau} \right), \quad (15)$$

579 while Equation (14) leads to

$$p(\mathbf{X} | \mathbf{t}, \gamma, m, g, h, I_p = 0) = \int_{\mu} \int_{\Lambda} p(\Lambda | g, h) p(\mu | m, \gamma, \Lambda) \prod_n p(x_n | \mu, \Lambda) d\mu, d\Lambda. \quad (16)$$

580 For coding equal preference for the indicator values $I_p = 1$ and $I_p = 0$ we use a flat prior and hence

581 $P(I_p = 1) = P(I_p = 0) = 0.5$. The marginal likelihoods in Equations (15) and (16) can subsequently

582 be converted to the posterior probability for differential regulation of protein abundance

583 $P(I_p \equiv 1 | \mathbf{t}, \gamma, m, g, h)$:

$$P(I_p = 1 | \mathbf{t}, \gamma, m, g, h) = \frac{p(\mathbf{X} | \mathbf{t}, \gamma, m, g, h, I_p = 1)}{p(\mathbf{X} | \mathbf{t}, \gamma, m, g, h, I_p = 1) + p(\mathbf{X} | \mathbf{t}, \gamma, m, g, h, I_p = 0)} \quad (17)$$

584
 585 Taking the posterior probability $P(I_p = 1 | \mathbf{t}, \gamma, m, g, h)$ in Equation (17) as measure of differential
 586 protein regulation is justified by the fact that Bayesian model selection has Occam's razor built in⁴⁴.
 587 Posterior probability values $P(I_p = 1 | \mathbf{t}, \gamma, m, g, h)$ which are larger than 0.5 will only be observed if
 588 the more complex model ($I_p = 1$) provides a substantially better fit of the data \mathbf{X} and \mathbf{t} than the simpler
 589 model ($I_p = 0$).

590 *Inferring differentially regulated protein complexes.* Inference of differential regulation of
 591 protein complexes assumes that the assignment of proteins to protein complexes is known. All
 592 subsequent derivations assume thus that the protein complex c is defined as a set of proteins $C_c =$
 593 $[P_1, P_2, \dots, P_C]$ of cardinality C . We assume furthermore that a set of retention profiles $X_c =$
 594 $[X_{P_1}, X_{P_2}, \dots, X_{P_C}]$ and a corresponding set of phenotype descriptions $t_c = [t_{P_1}, t_{P_2}, \dots, t_{P_C}]$ is available.
 595 If the retention profiles in X_c and thus the corresponding phenotype characteristics in t_c can at least in
 596 part be paired among all proteins which establish a complex, we have the subset $N_c = [n_1, n_2, \dots, n_K]$
 597 of samples for which complete observations are available. To prepare inferring differentially regulated
 598 protein complexes we may in this situation aggregate the protein specific retention profiles to a column
 599 concatenated matrix \mathbf{Y}_c which represents all retention profiles of the entire complex. Denoting the
 600 selection of the n^{th} row of matrix X_{P_c} as $X_{P_c}[n]$ and column wise row concatenation of row vectors
 601 $X_{P_c}[n]$ and $X_{P_{c+1}}[n]$ as $[X_{P_c}[n], [X_{P_{c+1}}[n]]$, we obtain

$$\mathbf{Y}_c = \begin{pmatrix} X_{P_1}[n_1], & X_{P_2}[n_1], & \dots, & X_{P_C}[n_1] \\ X_{P_1}[n_2], & X_{P_2}[n_2], & \dots, & X_{P_C}[n_2] \\ \vdots & \dots & \dots & \vdots \\ X_{P_1}[n_K], & X_{P_2}[n_K], & \dots, & X_{P_C}[n_K] \end{pmatrix},$$

603 as protein complex specific matrix of retention profiles and $\mathbf{u}_c = [t_{P_1}[n_1], t_{P_1}[n_2], \dots, t_{P_1}[n_K]]^T$ as
 604 protein complex specific phenotype vector. Inference of differentially regulated complexes is now a
 605 straightforward application of Equation (15), Equation (16) and Equation (17). We have just got to
 606 replace the protein specific retention profiles \mathbf{X} in these equations with the protein complex specific

607 retention profiles Y_c and exchange the phenotype characterization t with the phenotype
608 characterization of the protein complex u_c . Pooling of retention profiles requires in addition to the
609 assumptions which led to the DAG in Figure 1 no additional assumptions. While this is an advantage
610 of the approach, we have to consider that pooling of samples requires complete sets of paired retention
611 profiles which have to be available for all proteins which aggregate to the complex. In practice
612 measurement errors will lead to random dropouts and thus to a potentially small number of samples
613 where all data is available. To avoid such information loss by pairing of samples, we propose an
614 additional approach for assessing differentially regulated protein complexes by Bayesian model
615 probabilities. For assessing differential regulation of protein complex c we apply Equations (15), (16)
616 and (17) for every protein $P_k \in C_c$ separately. Following the assessment on protein level, differential
617 expression of protein complex c is coded via a binary indicator variable C_c . Assuming conditional
618 independence among proteins we may represent this proposition by the DAG in **Supplementary Fig.**
619 **S7b**. The DAG leads for the posterior probability of differential regulation of protein complex C_c
620 finally to Equation (18).

621

$$P(C_c \equiv 1 | \mathbf{X}_{P_1}, \dots, \mathbf{X}_{P_C}, \mathbf{t}_{P_1}, \dots, \mathbf{t}_{P_C}, \gamma, m, g, h) = \frac{P(C_c \equiv 1) \prod_{P_k \in C_c} p(\mathbf{X}_{P_k} | \mathbf{t}_{P_k}, \gamma, m, g, h, I_{P_k} \equiv 1)}{\sum_{I=0}^1 P(C_c \equiv I) \prod_{P_k \in C_c} p(\mathbf{X}_{P_k} | \mathbf{t}_{P_k}, \gamma, m, g, h, I_{P_k} \equiv I)}, \quad (18)$$

622

623 with $P_k \in C_c$ denoting all proteins which aggregate to the protein complex c . The prior probability for
624 protein complex c being differentially regulated, $P(C_c)$, is assumed to be identical for both indicator
625 values and thus $P(C_c = 1) = P(C_c = 0) = 0.5$. The expression $p(\mathbf{X}_{P_k} | \mathbf{t}_{P_k}, \gamma, m, g, h, I_{P_k} = [0,1])$ denotes
626 for $I = [0, 1]$ the marginal likelihoods we obtain with the model in **Supplementary Fig. S7a** for
627 protein P_k according to Supplementary Equations (15) and (16).

628

629 **Software implementation and data visualisation**

630 The command-line version of PCprophet was implemented and visualised using Python, together with
631 third-party packages including SciPy, Pandas⁴⁵, scikit-learn³⁷, NetworkX⁴⁶, and Matplotlib⁴⁷.

632

633 Source code availability

634 PCprophet is open-access and freely available for academic purposes at

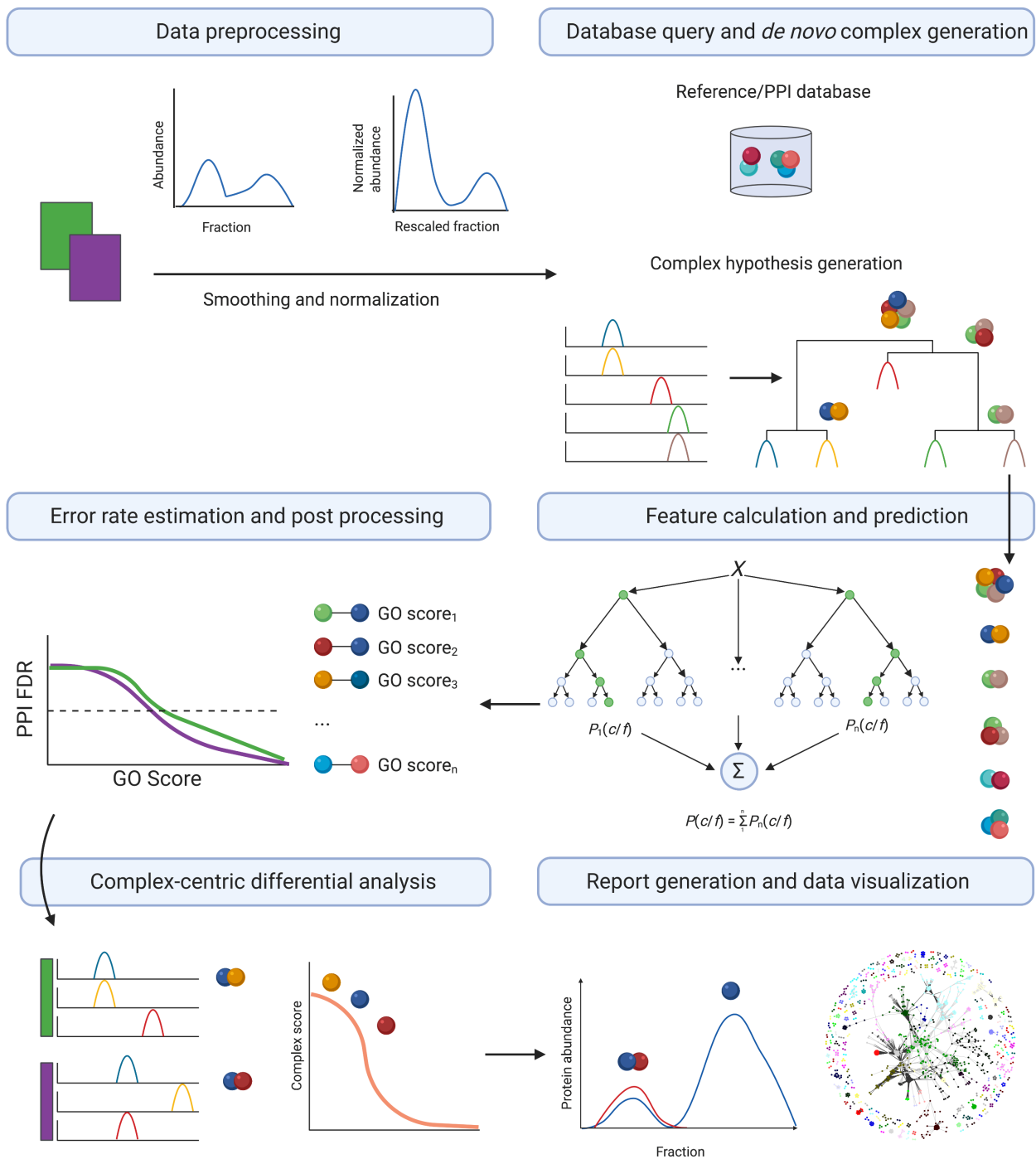
635 <https://github.com/fossatiA/PCprophet> under the MIT License.

636

637 References

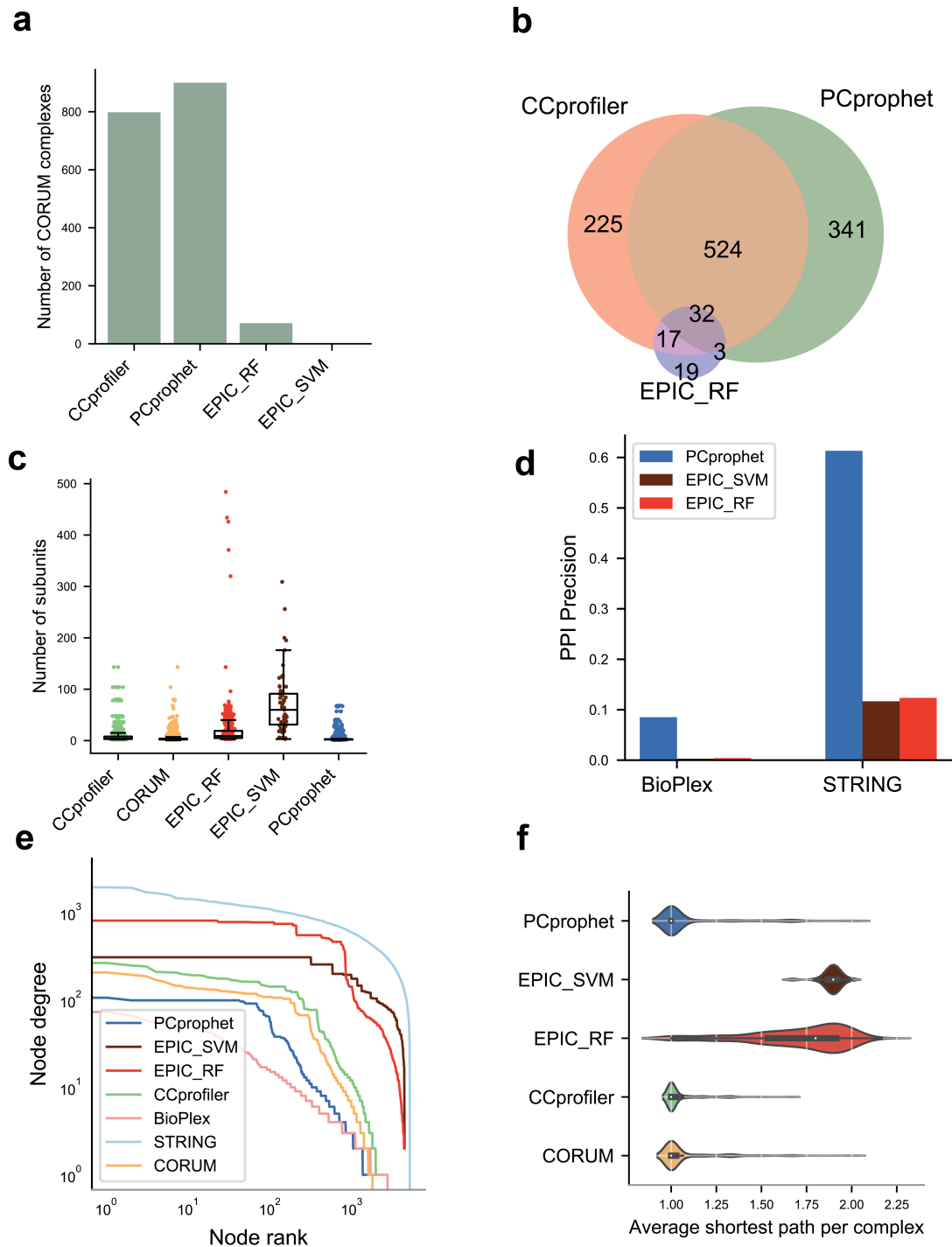
- 638 1. Marsh, J.A. & Teichmann, S.A. Structure, dynamics, assembly, and evolution of protein
639 complexes. *Annu Rev Biochem* **84**, 551-575 (2015).
- 640 2. Pan, J. et al. Interrogation of Mammalian Protein Complex Structure, Function, and
641 Membership Using Genome-Scale Fitness Screens. *Cell Syst* **6**, 555-568 e557 (2018).
- 642 3. Sowmya, G., Breen, E.J. & Ranganathan, S. Linking structural features of protein complexes
643 and biological function. *Protein Sci* **24**, 1486-1494 (2015).
- 644 4. Spirin, V. & Mirny, L.A. Protein complexes and functional modules in molecular networks.
645 *Proc Natl Acad Sci U S A* **100**, 12123-12128 (2003).
- 646 5. Salas, D., Stacey, R.G., Akinlaja, M. & Foster, L.J. Next-generation Interactomics:
647 Considerations for the Use of Co-elution to Measure Protein Interaction Networks. *Mol Cell*
648 *Proteomics* **19**, 1-10 (2020).
- 649 6. Scott, N.E. et al. Interactome disassembly during apoptosis occurs independent of caspase
650 cleavage. *Mol Syst Biol* **13**, 906 (2017).
- 651 7. Crozier, T.W.M., Tinti, M., Larance, M., Lamond, A.I. & Ferguson, M.A.J. Prediction of
652 Protein Complexes in *Trypanosoma brucei* by Protein Correlation Profiling Mass Spectrometry
653 and Machine Learning. *Mol Cell Proteomics* **16**, 2254-2267 (2017).
- 654 8. Hu, L.Z. et al. EPIC: software toolkit for elution profile-based inference of protein complexes.
655 *Nat Methods* **16**, 737-742 (2019).
- 656 9. Kirkwood, K.J., Ahmad, Y., Larance, M. & Lamond, A.I. Characterization of native protein
657 complexes and protein isoform variation using size-fractionation-based quantitative proteomics.
658 *Mol Cell Proteomics* **12**, 3851-3873 (2013).
- 659 10. Heusel, M. et al. Complex-centric proteome profiling by SEC-SWATH-MS. *Mol Syst Biol* **15**,
660 e8438 (2019).
- 661 11. McBride, Z. et al. A Label-free Mass Spectrometry Method to Predict Endogenous Protein
662 Complex Composition. *Mol Cell Proteomics* **18**, 1588-1606 (2019).
- 663 12. Stacey, R.G., Skinnider, M.A., Scott, N.E. & Foster, L.J. A rapid and accurate approach for
664 prediction of interactomes from co-elution data (PrInCE). *BMC Bioinformatics* **18**, 457 (2017).
- 665 13. Giurgiu, M. et al. CORUM: the comprehensive resource of mammalian protein complexes-
666 2019. *Nucleic Acids Res* **47**, D559-D563 (2019).
- 667 14. Szklarczyk, D. et al. STRING v11: protein-protein association networks with increased
668 coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic*
669 *Acids Res* **47**, D607-D613 (2019).
- 670 15. Huttlin, E.L. et al. Architecture of the human interactome defines protein communities and
671 disease networks. *Nature* **545**, 505-509 (2017).
- 672 16. Huttlin, E.L. et al. The BioPlex Network: A Systematic Exploration of the Human Interactome.
673 *Cell* **162**, 425-440 (2015).
- 674 17. Heusel, M. et al. A global screen for assembly state changes of the mitotic proteome by SEC-
675 SWATH-MS. *Cell Syst* (2019).
- 676 18. Kristensen, A.R., Gsponer, J. & Foster, L.J. A high-throughput approach for measuring
677 temporal changes in the interactome. *Nat Methods* **9**, 907-909 (2012).
- 678 19. Havugimana, P.C. et al. A census of human soluble protein complexes. *Cell* **150**, 1068-1081
679 (2012).

- 680 20. Oughtred, R. et al. The BioGRID interaction database: 2019 update. *Nucleic Acids Res* **47**,
681 D529-D541 (2019).
- 682 21. Hamazaki, J. et al. A novel proteasome interacting protein recruits the deubiquitinating enzyme
683 UCH37 to 26S proteasomes. *EMBO J* **25**, 4524-4536 (2006).
- 684 22. Lasker, K. et al. Molecular architecture of the 26S proteasome holocomplex determined by an
685 integrative approach. *Proc Natl Acad Sci U S A* **109**, 1380-1387 (2012).
- 686 23. Huang, D.T. et al. E2-RING expansion of the NEDD8 cascade confers specificity to cullin
687 modification. *Mol Cell* **33**, 483-495 (2009).
- 688 24. Kohroki, J., Nishiyama, T., Nakamura, T. & Masuho, Y. ASB proteins interact with Cullin5
689 and Rbx2 to form E3 ubiquitin ligase complexes. *FEBS Lett* **579**, 6796-6802 (2005).
- 690 25. Lowe, N. et al. Analysis of the expression patterns, subcellular localisations and interaction
691 partners of Drosophila proteins using a pigP protein trap library. *Development* **141**, 3994-4005
692 (2014).
- 693 26. Collins, M.O. et al. Molecular characterization and comparison of the components and
694 multiprotein complexes in the postsynaptic proteome. *J Neurochem* **97 Suppl 1**, 16-23 (2006).
- 695 27. Antonyamy, S. et al. Crystal structure of the human PRMT5:MEP50 complex. *Proc Natl Acad*
696 *Sci U S A* **109**, 17960-17965 (2012).
- 697 28. Scoumanne, A., Zhang, J. & Chen, X. PRMT5 is required for cell-cycle progression and p53
698 tumor suppressor function. *Nucleic Acids Res* **37**, 4965-4976 (2009).
- 699 29. Gu, Z. et al. The p44/wdr77-dependent cellular proliferation process during lung development
700 is reactivated in lung cancer. *Oncogene* **32**, 1888-1900 (2013).
- 701 30. Rost, H.L. et al. TRIC: an automated alignment strategy for reproducible protein quantification
702 in targeted proteomics. *Nat Methods* **13**, 777-783 (2016).
- 703 31. R., Q. C4.5: Programs for Machine Learning. (Morgan Kaufmann Publishers, 1993).
- 704 32. Breiman, L. Random forests. *Mach Learn* **45**, 5-32 (2001).
- 705 33. John, G.H. & Langley, P. in Eleventh Conference on Uncertainty in Artificial Intelligence 338-
706 345 (Morgan Kaufmann, 1995).
- 707 34. Cortes, C. & Vapnik, V. Support-Vector Networks. *Mach Learn* **20**, 273-297 (1995).
- 708 35. Lecessie, S. & Vanhouwelingen, J.C. Ridge Estimators in Logistic-Regression. *Appl Stat-J Roy*
709 *St C* **41**, 191-201 (1992).
- 710 36. Vert, J.P., Tsuda, K. & Schoelkopf, B. in Kernel Methods in Computational Biology 35-70
711 (MIT Press, Cambridge, MA, USA; 2004).
- 712 37. Pedregosa, F. et al. Scikit-learn: Machine Learning in Python. *J Mach Learn Res* **12**, 2825-2830
713 (2011).
- 714 38. E., F., M.A., H. & I.H., W. The WEKA Workbench. Online Appendix for "Data Mining:
715 Practical Machine Learning Tools and Techniques", Edn. Fourth Edition. (Morgan Kaufmann,
716 2016).
- 717 39. Ashburner, M. et al. Gene ontology: tool for the unification of biology. The Gene Ontology
718 Consortium. *Nat Genet* **25**, 25-29 (2000).
- 719 40. Carbon, S. et al. AmiGO: online access to ontology and annotation data. *Bioinformatics* **25**,
720 288-289 (2009).
- 721 41. The Gene Ontology, C. The Gene Ontology Resource: 20 years and still GOing strong. *Nucleic*
722 *Acids Res* **47**, D330-D338 (2019).
- 723 42. Wang, J.Z., Du, Z., Payattakool, R., Yu, P.S. & Chen, C.F. A new method to measure the
724 semantic similarity of GO terms. *Bioinformatics* **23**, 1274-1281 (2007).
- 725 43. UniProt, C. UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res* **47**, D506-D515
726 (2019).
- 727 44. Bernardo, J.M. & Smith, A.F.M. Bayesian Theory. (John Wiley & Sons, Inc., 1994).
- 728 45. W., M. in The 9th Python in Science Conference 51-56 (2010).
- 729 46. A.A., H., D.A., S. & P.J., S. in The 7th Python in Science Conference (SciPy2008) (Pasadena,
730 CA USA; 2008).
- 731 47. Hunter, J.D. Matplotlib: A 2D graphics environment. *Comput Sci Eng* **9**, 90-95 (2007).
- 732

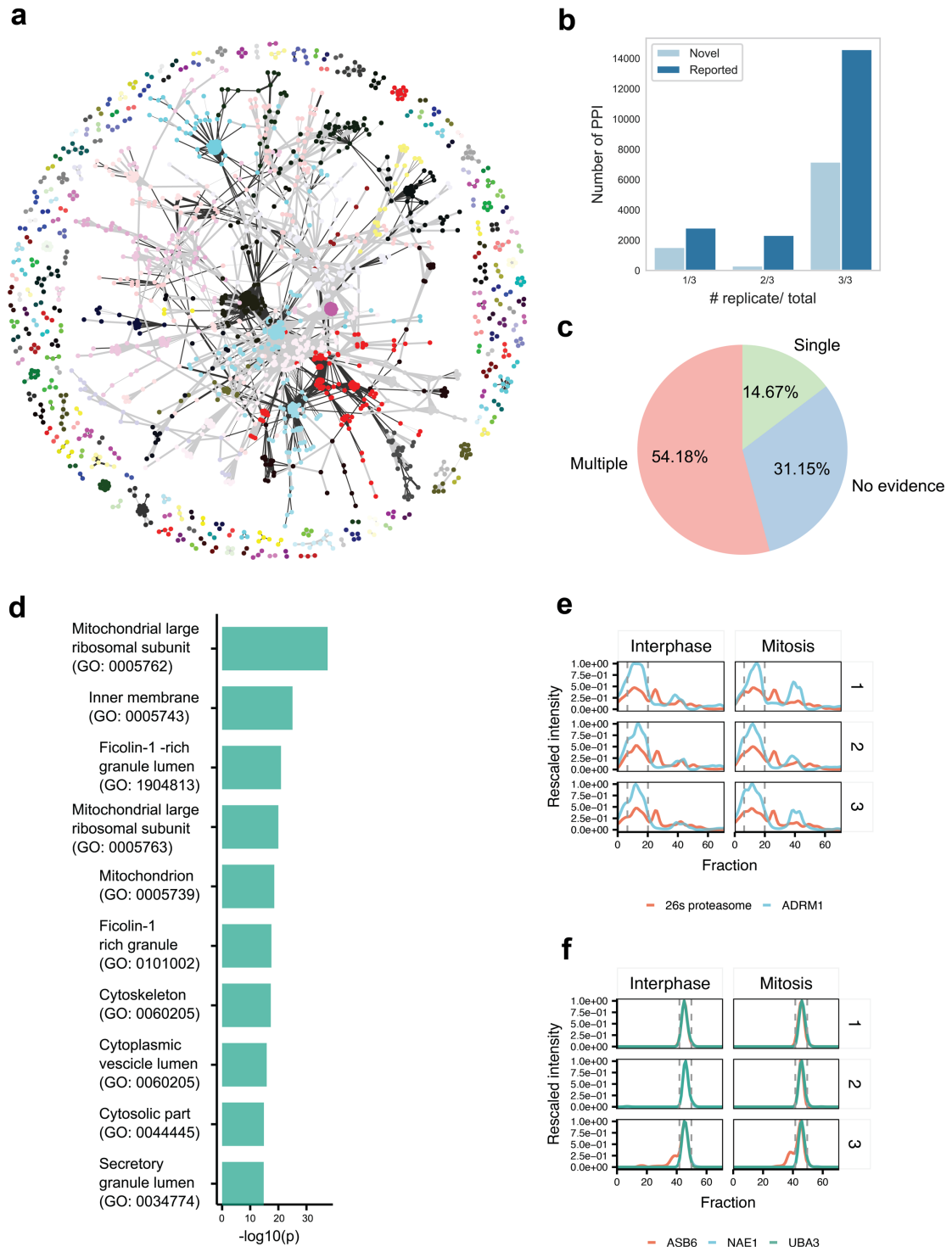


733
 734 **Fig. 1.** The framework of PCprophet. It consists of the six major modules including (i) data pre-
 735 processing, (ii) database query and *de novo* complex (i.e. hypothesis) generation, (iii) feature
 736 calculation and prediction, (iv) error estimation and post-prediction processing, (v) complex-centric
 737 differential analysis, and (vi) report generation and data visualisation.

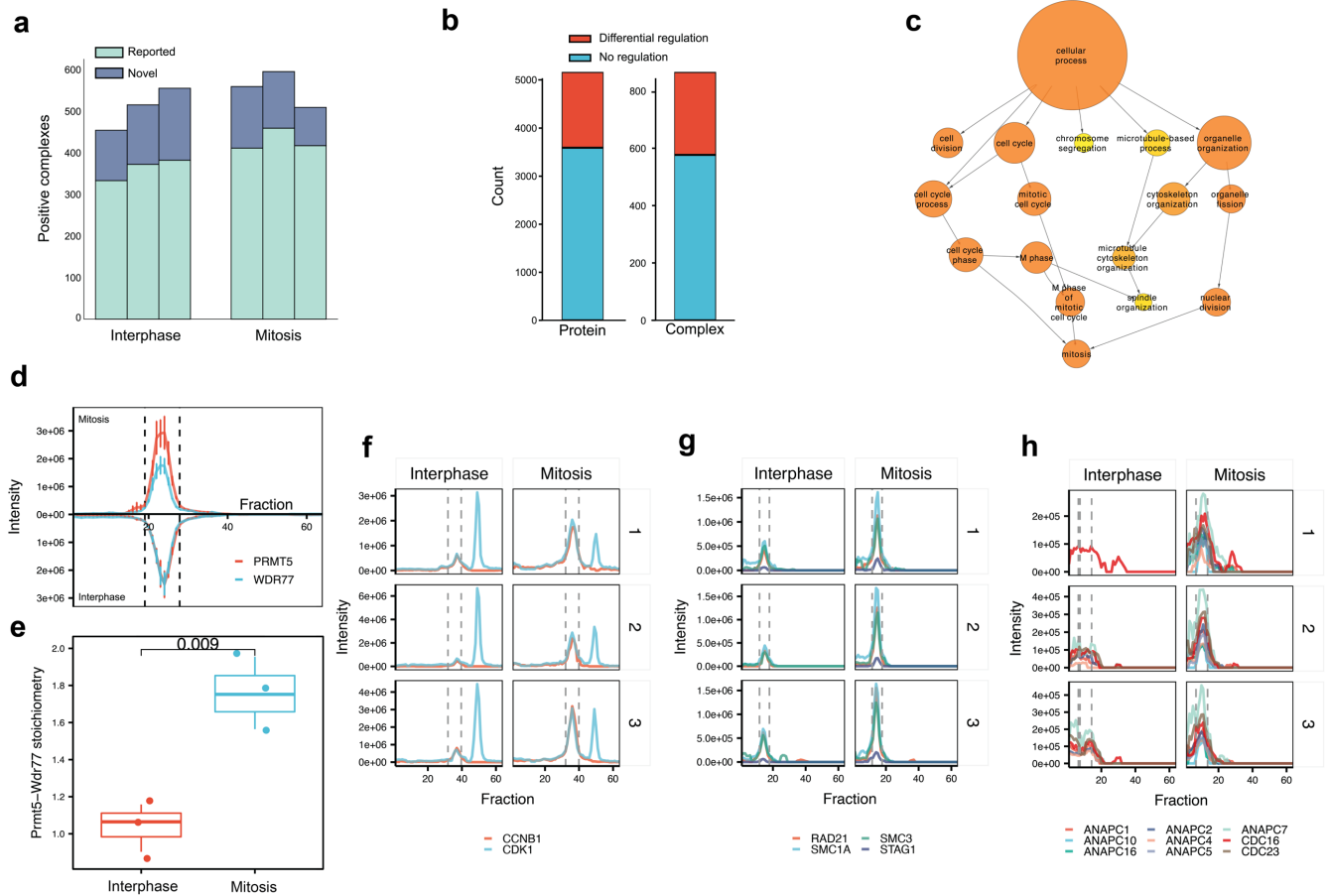
738



739
 740 **Fig. 2.** Benchmarking PCprophet against existing tools for protein complex profiling and prediction. **a**,
 741 The numbers of CORUM complexes recovered and the numbers of overlapping complexes by the
 742 assessed tools. **b**, Absolute number of CORUM complexes recovered by each tool. **c**, Number of
 743 subunits per complex predicted and identified by different tools. Boxplot shows the medians and the
 744 ticks represent standard deviation. **d**, The precision values (refer to the ‘Methods’ section for more
 745 details) of PPI prediction for *de novo* protein complex prediction tools. **e**, Log-log plot showing the
 746 degree distribution of the network generated by each tool *versus* ground-truth databases (STRING,
 747 BioPlex and CORUM). **f**, Distribution of shortest path per complex across all subunits, as reported by
 748 the indicated tools. The medians are highlighted in white dots.



749
 750 **Fig. 3.** Evaluation of *de novo* prediction using PCprophet. **a**, The PPI map generated by PCprophet
 751 from HeLa cell proteomic data. Edge width represents the number of technical replicates for which a
 752 particular PPI was found. Black edge are novel PPIs and grey edges are reported PPIs. Protein
 753 communities are highlighted in different node colours. **b**, Number of novel and reported PPIs across
 754 all conditions within technical replicates (i.e. interphase and mitosis). **c**, Annotation of novel PPIs (i.e.
 755 not documented in the CORUM databases) in PPI databases (STRING, BioPlex, BioGrid). **d**,
 756 Enrichment analysis for GO Cellular Component for PPIs without prior evidence in any database. **e**,
 757 26S proteasome and ADRM1 coelution from interphase and mitotic cells. **f**, Co-elution of ASB6 with
 758 the NAE1 and UBA3 complex.



759
 760 **Fig. 4.** Differential analysis of complexes across the cell cycle states tested. **a**, Absolute number of
 761 novel and reported complexes in the indicated cell cycle stages. **b**, Stacked histogram for differentially
 762 regulated proteins (n=1518) and differentially regulated complexes (n=238). **c**, Enrichment for GO
 763 Biological Process for differentially regulated proteins between the two conditions using as background
 764 all proteins identified. Node size represents number of proteins within the particular category. Nodes
 765 colour represents Bonferroni adjusted p value, ranging from $p=10^{-3}$ (yellow) to $p=10^{-8}$ (orange). **d**,
 766 Mirror plot for co-elution profiles of Prmt5-Wdr77 complex in mitosis (upper positive Y axis) and
 767 interphase (negative y axis). Values were averaged across the three replicates for each condition and
 768 bar represents standard error of the mean. **e**, Bar-plot of Prmt5/Wdr77 complex stoichiometry in
 769 interphase (red, mean=1.04) and mitosis (blue, mean=1.75). **f**, Co-elution profile for the Ccnb1/CKD1
 770 complex **(G)** Co-elution profile for the Anaphase promoting complex. **h**, Co-elution profile for the
 771 cohesion complex.

772 **Supplementary Methods**

773 **Mathematical details for Bayesian inference of differential regulation of protein abundance and** 774 **protein complexes**

775 We now present further mathematical details how we may express the marginal likelihoods in
776 Equations (15) and (16). We start this derivation by expressing the prior densities over Λ and μ for the
777 simpler model $I_p = 0$ and the densities over $\mu_\tau \forall \tau$ for the more complex model $I_p = 1$. As is mentioned
778 above, the prior over the precision matrix Λ is coded as a product of Gamma densities. With $\Lambda =$
779 $diag([\lambda_1, \dots, \lambda_D])$ and D denoting the input dimension (number of columns) of X we get

$$p(\Lambda|g, h) = \prod_{d=1}^D \left(\frac{h^g}{\Gamma(g)} \lambda_d^{g-1} \exp(-h\lambda_d) \right), \quad (19)$$

780
781 where $\Gamma(g) = \int_{x=0}^{\infty} x^{g-1} \exp(-x) dx$ denotes the Gamma function. The multivariate Gaussian prior
782 over μ for the non-differentially regulated case $I_p = 0$ is

$$p(\mu|\Lambda, \gamma, m) = (2\pi)^{-\frac{D}{2}} \gamma^{\frac{D}{2}} \|\Lambda\|^{-\frac{1}{2}} \exp(-0.5\gamma(\mu - m)^T \Lambda(\mu - m)), \quad (20)$$

783
784 where $\|\Lambda\|$ denotes the determinant of the precision matrix. Finally, we get the multivariate Gaussian
785 prior over μ for the differentially regulated case $I_p = 1$ as

$$p(\mu|\Lambda, \gamma, m) = \prod_{\tau=1}^{\#t} \left((2\pi)^{-\frac{D}{2}} \gamma^{\frac{D}{2}} \|\Lambda\|^{-\frac{1}{2}} \exp(-0.5\gamma(\mu_\tau - m)^T \Lambda(\mu_\tau - m)) \right). \quad (21)$$

786
787 To calculate the marginal likelihood for the model $I_p = 1$ we integrate Equation (13) first with respect
788 to all μ_τ and then with respect to Λ to get

$$p(\mathbf{X}|\mathbf{t}, \gamma, m, g, h, I_p \equiv 1) = \left(\frac{h^g}{\Gamma(g)} \right)^D (2\pi)^{-\frac{D \cdot N}{2}} \gamma^{\frac{D \cdot \#t}{2}} \prod_{\tau=1}^{\#t} (\gamma + n_\tau)^{-\frac{D}{2}} \prod_{d=1}^D \Gamma(\hat{g}) / \hat{h}_d^{\hat{g}}, \text{ where}$$

$$\hat{g} = g + \frac{N}{2}, \quad (22)$$

$$\hat{h}_d = h + 0.5 \left(\mathbf{t} \gamma m_d^2 + \sum_{n=1}^N x_{n,d}^2 - \sum_{\tau=1}^{\#t} (\xi_d^\tau)^2 / \gamma + n_\tau \right),$$

$$\xi^\tau = \gamma m + \sum_{n|t_n=\tau} x_n,$$

789

790 and we use N to denote the number of all samples and n_τ to denote the number of samples which have
791 phenotype level τ . The final expression for the marginal likelihood of the simpler model $I_p = 0$ from
792 Equation (16) is easily obtained from Equation (22). We just have to replace n_τ with N and $\#t$ with 1.
793 We should note that for numerical stability we calculate $\log(p(X|t, \gamma, m, g, h, I_p))$. The calculations
794 reported in this paper set the hyper parameters for g , h and γ to $g = 0.8$, $h = 1.5$ and $\gamma = 0.025$. As prior
795 location m we use the sample mean or set $m = 0$.
796

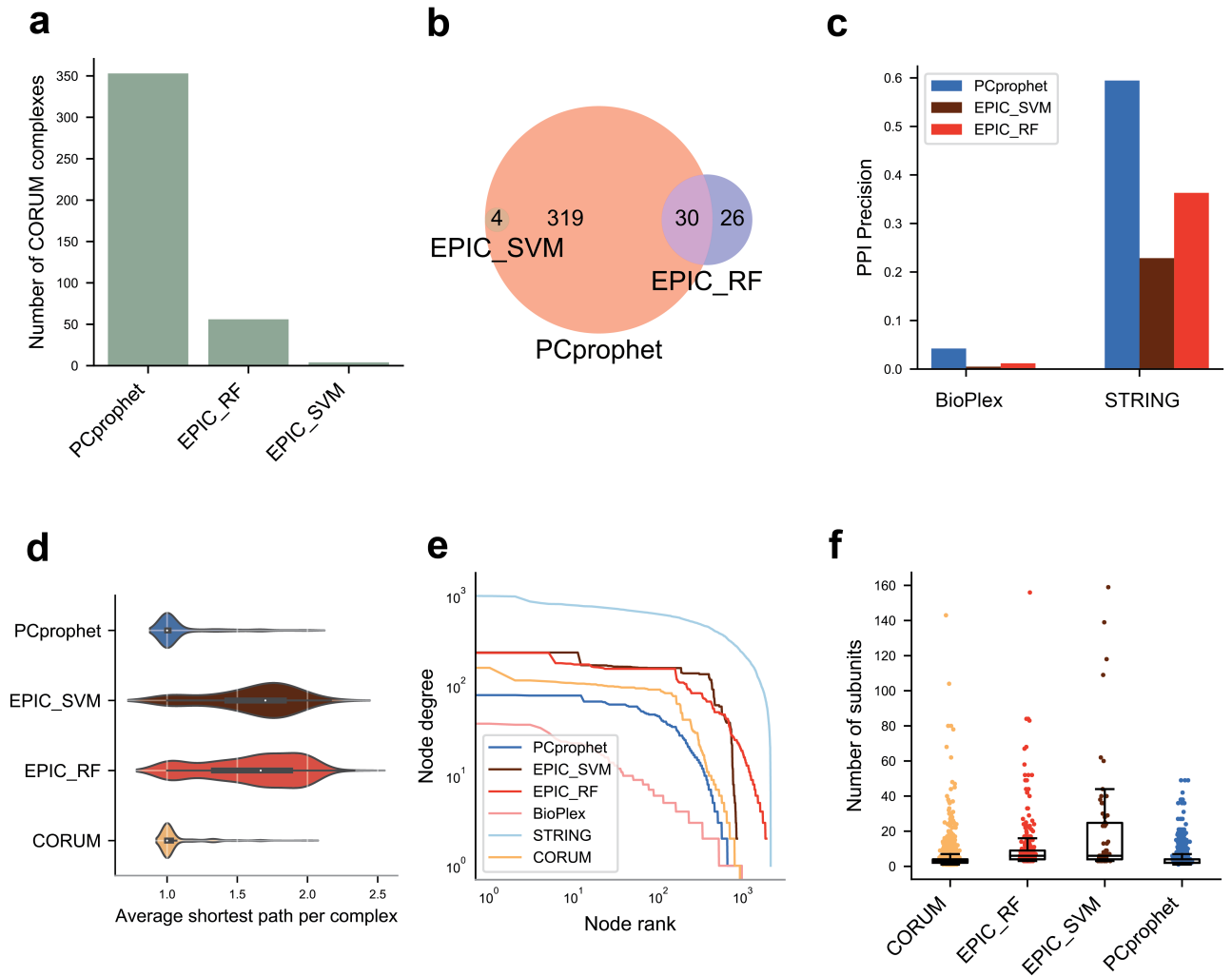
797 **Supplementary Results**

798 **Optimization of PCprophet machine learning framework for complex prediction**

799 In order to reach optimal performance in correctly classifying protein complex signals from the co-
800 fractionation datasets, we explored different types of machine learning strategies and their performance
801 to recall a set of manually curated protein complex signals in a previously published dataset¹⁰. To train
802 the machine learning models, we used manually annotated data (refer to the ‘**Methods**’ section for
803 more details) using criteria similar to the strategy applied in Heusel *et al*¹⁰. As the negative complexes
804 significantly outnumbered the positive complexes (i.e. 738 vs. 242) based on our manual annotation,
805 we evaluated the performance of PCprophet on two instances of the input dataset: one where all the
806 negatives were used and the other where an equal number of negatives as positives were randomly
807 selected. We tested the performance of PCprophet based on five well-established machine learning
808 models using five-fold cross-validation including Decision Tree (J48)³¹, Random Forest³² (RF), Naïve
809 Bayes³³ (NB), Support Vector Machines³⁴ (SVM) and Logistic Regression³⁵ (LR) algorithms [Figure
810 comparing the performance of the different algorithms. We determined that the RF achieved its best
811 performance when the number of trees was set to 500 via a separate stratified 10-fold cross-validation
812 on the entire dataset (**Supplementary Fig. S8**); for all other algorithms, we used default parameters. In
813 the latter dataset, we performed 100 trials for this random selection procedure (**Supplementary Table**
814 **S4**). RF outperformed all the other machine-learning algorithms and selected as the algorithm for
815 PCprophet, for example, achieving an AUC of 0.991, accuracy of 96.9% and MCC of 0.916,
816 irrespective of the number of negative complexes used (**Supplementary Fig. 9; Supplementary Table**
817 **S4**). To build a balanced and unbiased classifier, we rebuilt the RF model on the dataset with equal
818 positive and negative complexes, on which RF achieved the best performance according to the 100
819 trials of 5-fold cross-validation. This rebuilt RF model is then used as the core predictor for PCprophet.

820 **Supplementary Figures**

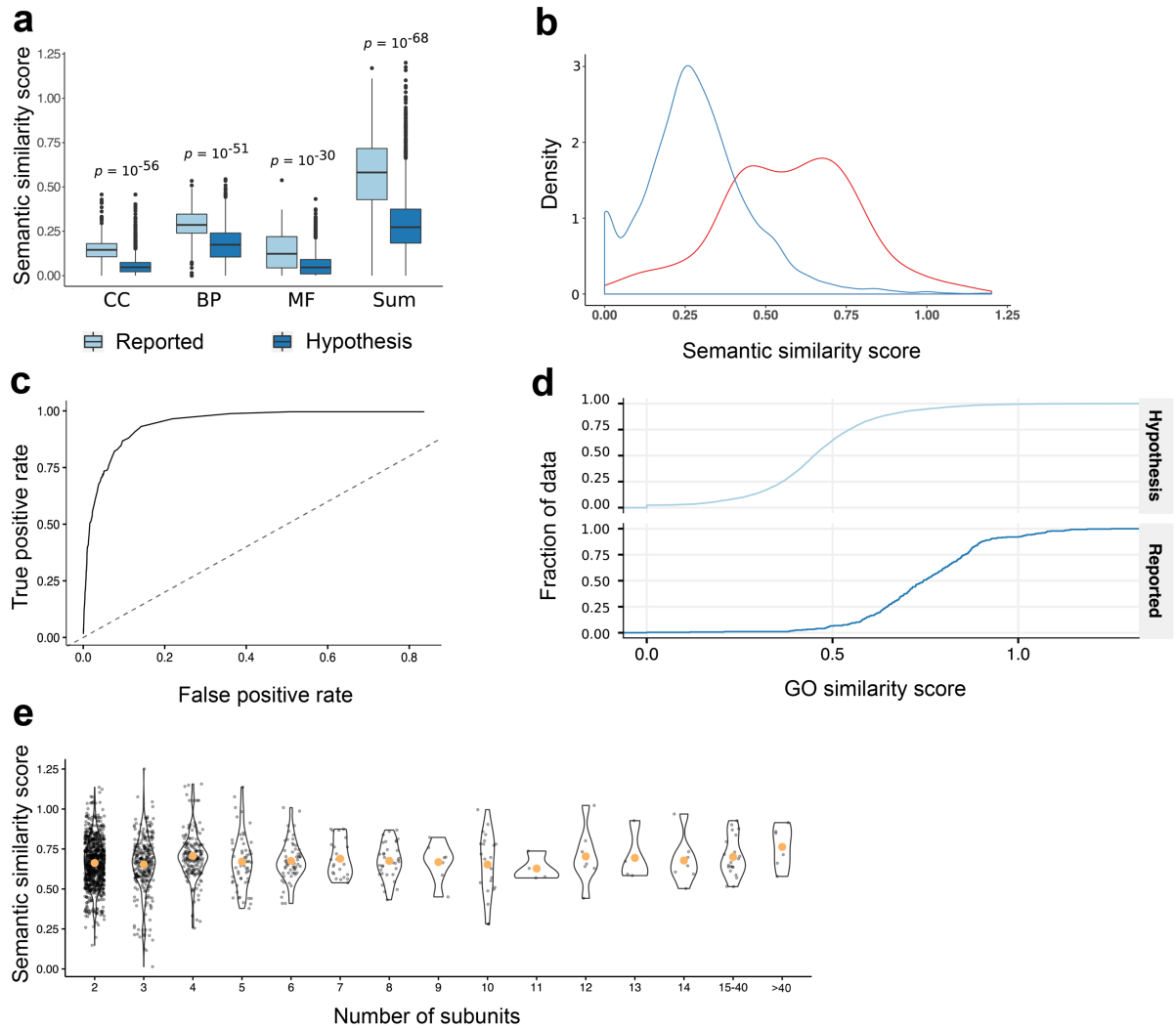
821



822

823 **Supplementary Fig. S1.** Benchmarking PCprophet with state-of-the-art software for complex profiling
824 on the DDA-SILAC dataset. **a**, Absolute number of CORUM complexes recovered by each tool. **b**,
825 Complex IDs overlap across all tools. **c**, Precision of PPI prediction for *de novo* protein complex
826 prediction tools. **d**, Distribution of shortest path per complex across all subunits. The medians are
827 highlighted using the white dots. **e**, Log-log plot showing the topology of the network generated by
828 each tool *versus* ground-truth databases. **f**, Number of subunits per complex across different tools.
829 Boxplot shows the medians and the ticks represent standard deviation.

830



831

832 **Supplementary Fig. S2.** Evaluation of GO score for estimating the false discovery rate. **a**, The boxplot

833 illustrating the separation of hypothesis and CORUM using the three individual ontologies or the

834 combination of the three (i.e. Molecular Function – MF; Biological Process – BP; Cellular Component

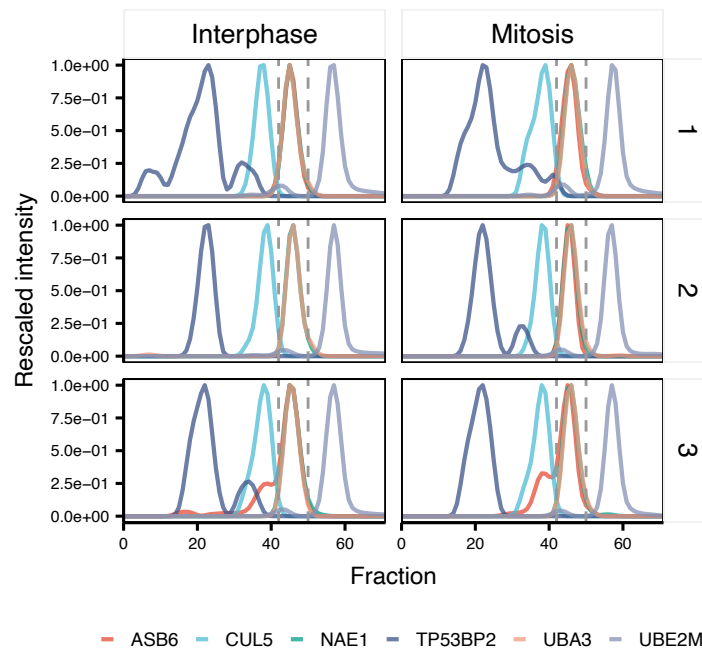
835 – CC). **b**, The density plot showing the separation of hypothesis and ground-truth CORUM database

836 using the sum of the three ontologies. **c**, Performance of GO term for separation of true and false PPIs.

837 **d**, Empirical cumulative distribution plot between hypothesis and reported complexes from CORUM.

838 **e**, Distribution of the sum of the three GO ontologies across different subunits size for reported

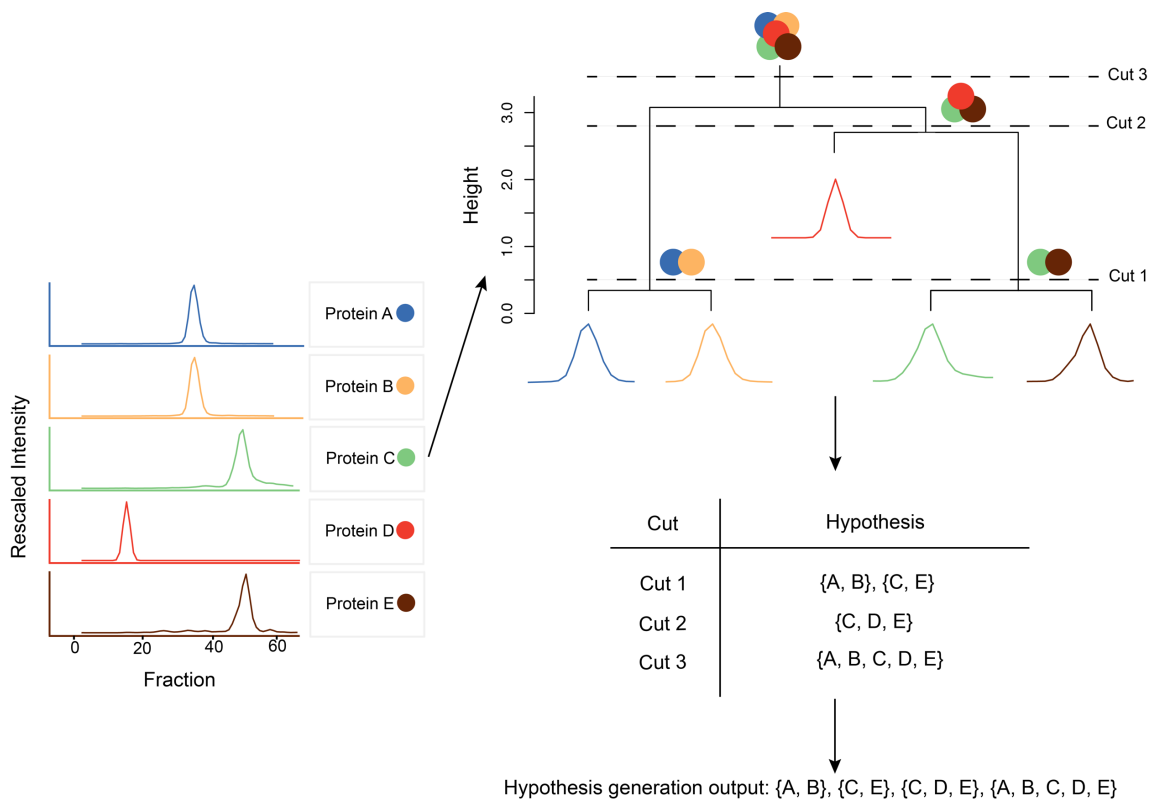
839 complexes.



840
841 **Supplementary Fig. S3.** Elution profiles of reported binders for ABS6 (CUL5) and NAE1-UBA3
842 (TP53BP2, UBE2M) across the two experimental conditions and the three replicates. The region
843 between the dotted lines represent the peak position of the novel ASB6-UBA3-NAE1 complex. The
844 absence of coelution of reported interactors between the dotted lines suggests the presence of a novel
845 complex rather than co-occurrence of complexes of similar size.

846
847

848

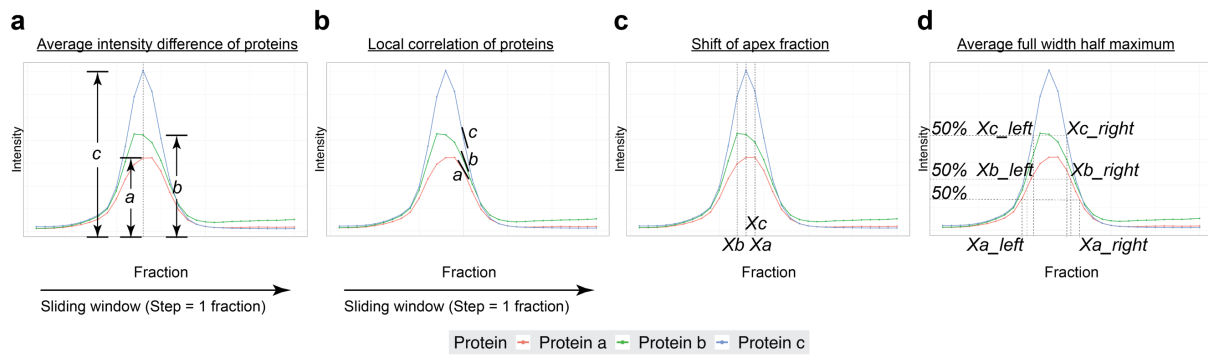


849

850 **Supplementary Fig. S4.** A conceptual illustration of hypothesis generation. Every protein is clustered
851 into possible complexes using Euclidian distance clustering. Following construction of a dendrogram,
852 all resulting clusters are retrieved by cutting at all heights (i.e. distances). This generates a
853 comprehensive set of all possible complexes in the data.

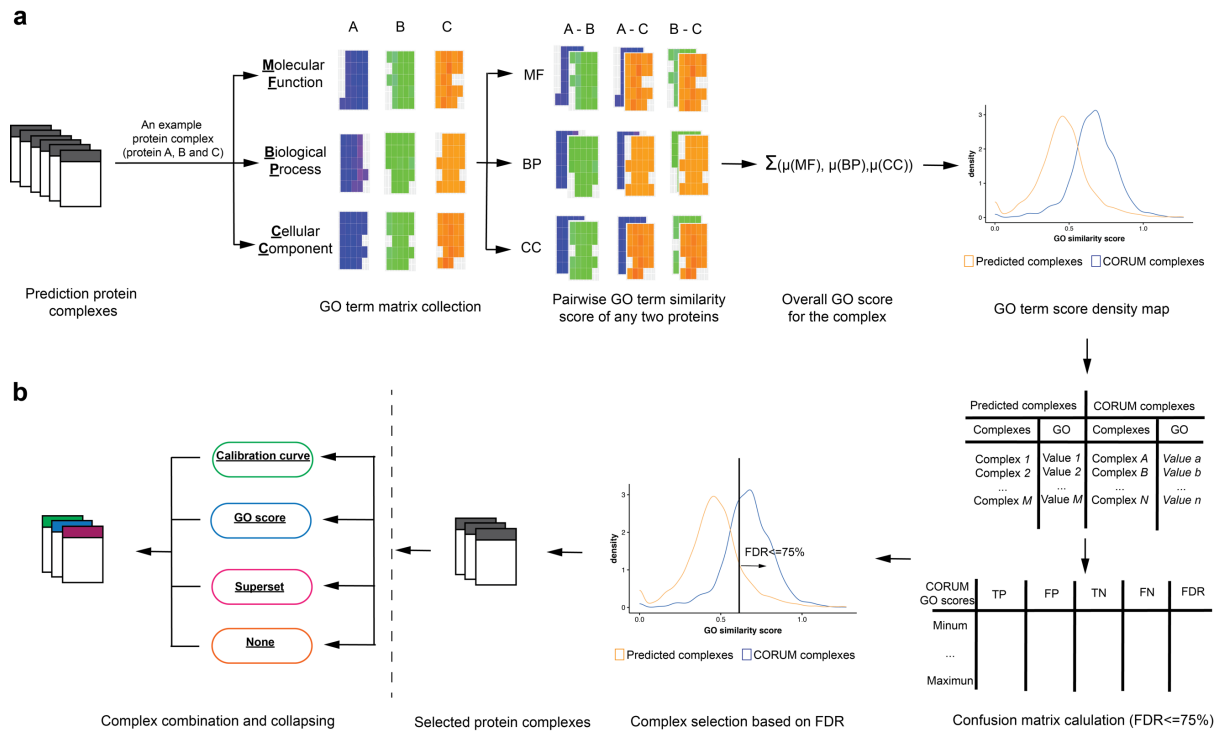
854

855



856

857 **Supplementary Fig. S5.** Feature calculation based on protein co-elution profiles. **a**, average intensity
858 difference of proteins. **b**, local correlation of proteins at each fraction. **c**, shift of apex fraction. **d**,
859 average full width half maximum.



860

861 **Supplementary Fig. S6.** Graphical illustration of post-prediction processing. **a**, GO term score filtering.

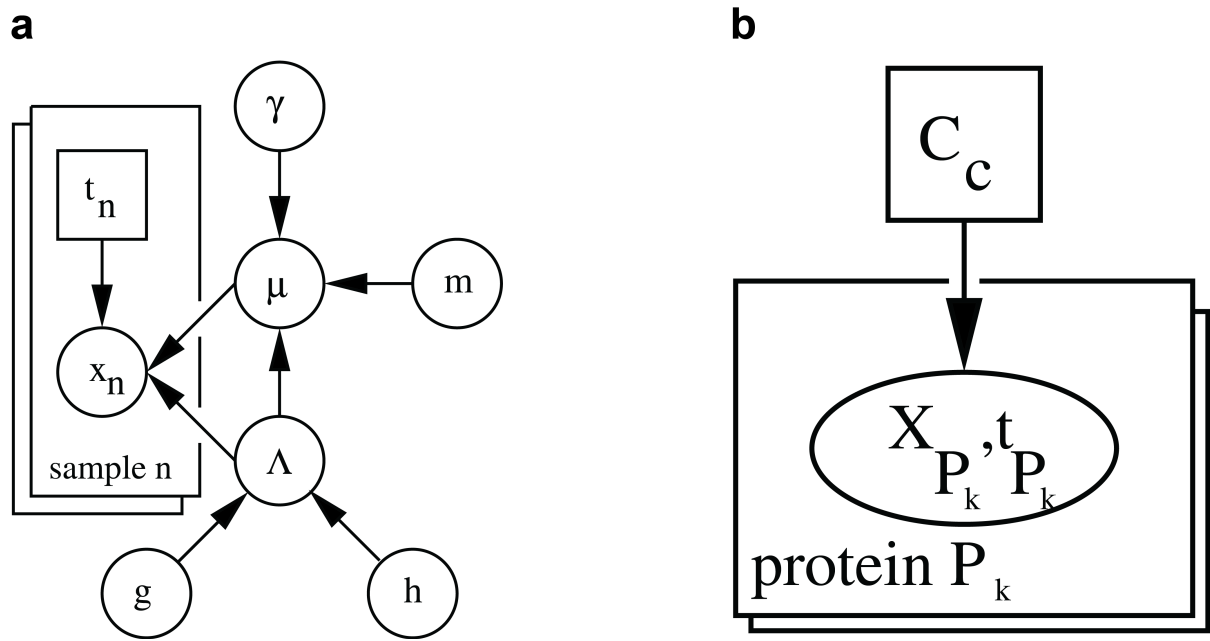
862 **b**, complex combination and collapsing. Positively predicted complexes either from the provided

863 database or PCprophet are decomposed into PPI and pairwise metrics are calculated based on semantic

864 similarity between the different ontologies which is then filtered based on a user-defined FDR threshold.

865 Overlapping complexes are combined based on the user-defined criteria.

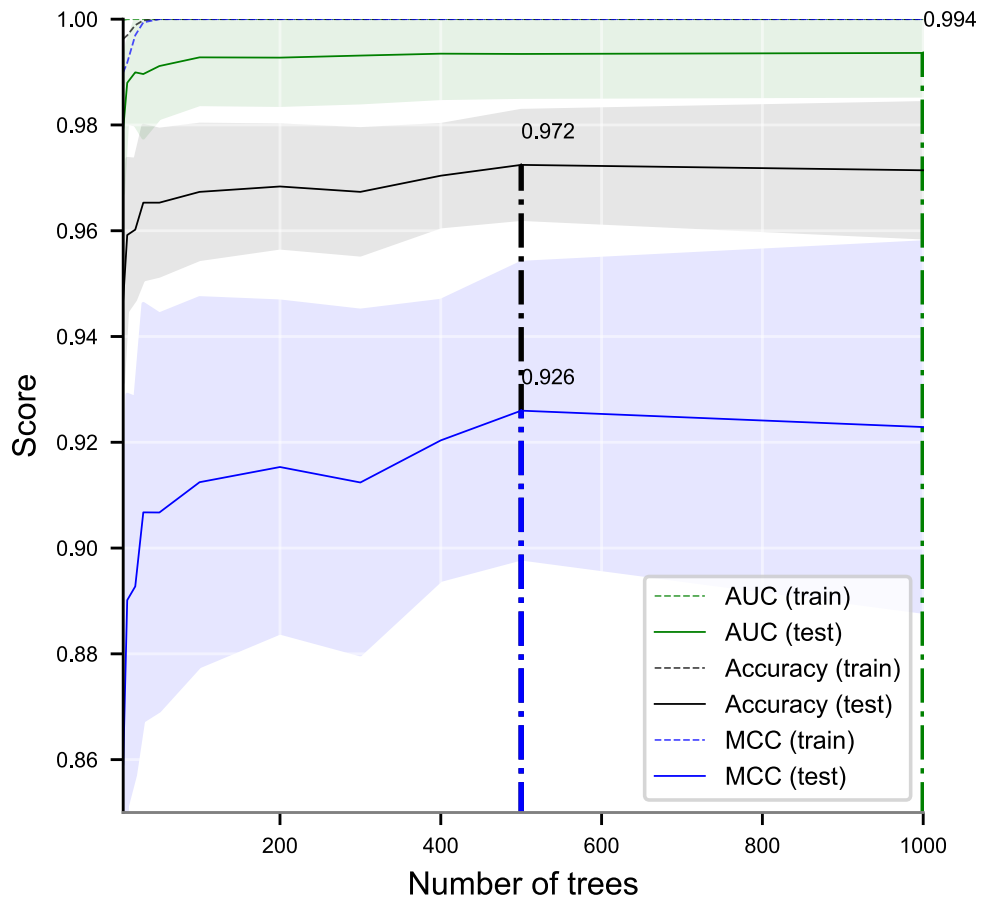
866



867
868

869 **Supplementary Figure S7.** Using Bayesian inference to analyse the difference of protein complexes
870 across conditions. **a**, An analytically tractable model for inferring differentially regulated proteins.
871 Variable μ denotes the mean of a multivariate Gaussian distribution over the protein abundance vector
872 x_n which depends on the phenotype state t_n . The prior over μ is a Gaussian distribution and
873 parameterized by m (the prior location), Λ and γ (together specifying the precision of the Gaussian
874 distribution). Variable Λ acts also as a precision matrix in the Gaussian $p(x_n|\mu, \Lambda)$. The prior over Λ
875 is a diagonal Wishart distribution (a product of Gamma distributions) which is parameterised by the hyper
876 parameters g and h . The conditional distribution $p(\mu, \Lambda|g, h, m, \gamma)$ which is represented by this DAG is
877 referred to as Normal-Wishart distribution and allows for an analytical calculation of Bayes factors. **b**,
878 A probabilistic model for inferring differentially regulated protein complexes. Variable C_c denotes the
879 state of differential regulation of a protein complex as binary variable. The elliptic node X_{P_k}, t_{P_k}
880 represents the marginal likelihood of the protein retention profiles X_{P_k} of protein $P_k \in C_c$ in
881 dependence of the phenotype characterization t_{P_k} as they arise from Equation (15) for $C_c = 1$ and from
882 Equation (16) for $C_c = 0$. For calculating the state of differential regulation of protein complex C_c we
883 make thus a conditional independence assumption among all contributing marginal likelihoods.
884

885



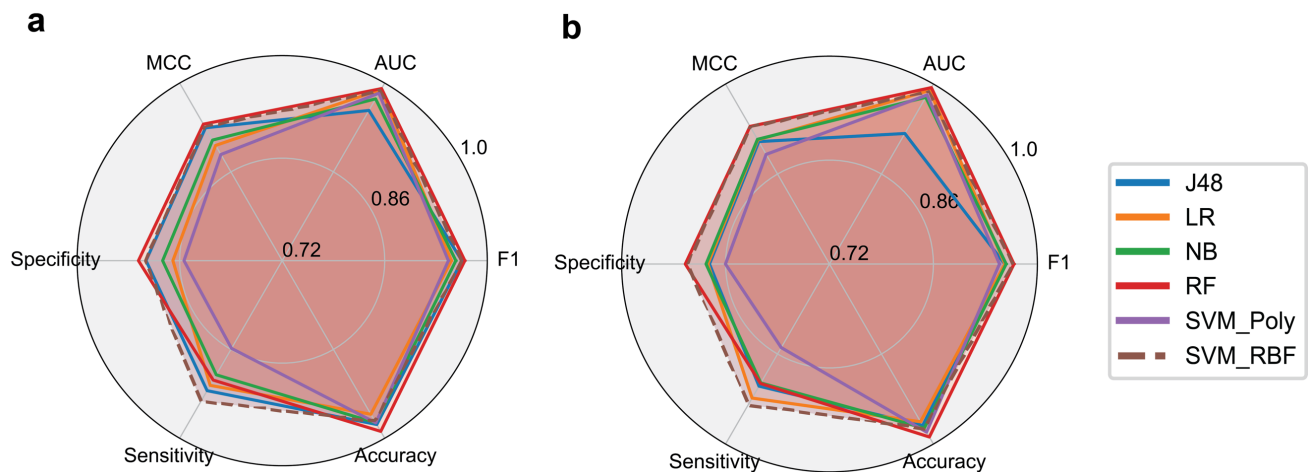
886
887 **Supplementary Fig. S8.** Determination of the number of trees in the RF model using three
888 performance evaluation measures, including accuracy, AUC and MCC, via stratified 10-fold cross-
889 validation on the entire dataset.

890

891

892

893



894

895 **Supplementary Fig. S9.** Radar plots demonstrating the prediction performance of PCprophet

896 via five-fold cross-validation. Performance parameters are AUC, accuracy, F1, MCC,

897 sensitivity and specificity. The analysis was based on the HEK293 dataset¹⁰ using **a**, equal

898 numbers of positives and negatives, which were randomly selected 100 times and **b**, positives

899 and all negatives. Coloured lines show the performance measures of different machine-learning

900 algorithms namely J48 decision tree (J48), linear regression (LR), Naïve Bayes (NB), Random

901 Forest (RF) and Support Vector Machine with either polynomial kernel (SVM_Poly) or radial

902 basis function kernel (SVM_RBF)

903

904 **Supplementary Tables**

905 **Supplementary Table S1.** AUC values evaluating the similarity between ground-truth
906 networks and prediction

907

908

AUC	Software	Δ AUC
PCprophet	0.096	0.013
EPIC_SVM	0.390	0.282
EPIC_RF	0.175	0.067
CCprofiler	0.122	0.014
BioPlex	0.040	0.068
STRING	0.137	0.029
CORUM	0.109	0

909

910 **Supplementary Table S2.** Detailed descriptions of the datasets applied for training and
911 evaluating PCprophet

Dataset	Cell line	Species	Data acquisition method	Separation technique
HEK293¹⁰	HEK293	<i>Homo sapiens</i>	SWATH	SEC (Size Exclusion Chromatography)
HeLa mitosis and interphase¹⁷	HeLa CCL2	<i>Homo sapiens</i>		
DDA-SILAC HeLa¹⁸	HeLa	<i>Homo sapiens</i>	DDA-SILAC	

912

913 **Supplementary Table S3.** An example demonstrating the complex collapsing step using three

914 protein complexes

Complex	Subunit	GO score	Calibration MW	Apparent MW
PC1	A, B	0.5	150,000	100,000
PC2	A, B, C	0.9	150,000	130,000
PC3	A, B, C, D	0.7	150,000	160,000

915

916 **Supplementary Table S4.** Prediction performance of PCprophet on manually annotated HEK293

917 datasets via 5-fold cross-validation using different numbers of negatives

	AUC	Accuracy	F1	MCC	Specificity	Sensitivity
	Using equal numbers of positives and negatives					
J48	0.957	96.531%	0.929	0.906	0.978	0.925
LR	0.990	95.102%	0.902	0.869	0.962	0.917
NB	0.975	95.714%	0.910	0.883	0.976	0.900
RF	0.991	96.939%	0.935	0.916	0.989	0.908
SVM_POLY	0.984	94.694%	0.887	0.854	0.976	0.858
SVM_RBF	0.988	96.531%	0.931	0.907	0.973	0.942
	Using positives and all negatives					
J48	0.923	95.603%	0.910	0.882	0.971	0.910
LR	0.990	95.622%	0.913	0.884	0.965	0.929
NB	0.979	95.796%	0.914	0.886	0.975	0.905
RF	0.994	97.268%	0.934	0.926	0.989	0.906
SVM_POLY	0.983	94.896%	0.891	0.860	0.981	0.850
SVM_RBF	0.988	96.684%	0.933	0.911	0.976	0.939

918

919