



39 expression control. The precise extent to which protein levels depend on mRNA abundances  
40 is still debated, and likely differs between genes and test systems<sup>21–23</sup>. However, some  
41 fundamental differences between mRNA and protein expression control have recently  
42 emerged. For example, many genes have coexpressed mRNAs due to their chromosomal  
43 proximity rather than any functional similarity<sup>19,24–26</sup>. Such non-functional mRNA coexpression  
44 results from stochastic transitions between active and inactive chromatin that affect wide  
45 genomic loci<sup>24,25,27</sup>, and transcriptional interference between closeby genes<sup>25,28</sup>. Importantly,  
46 coexpression of spatially close, but functionally unrelated genes is buffered at the protein  
47 level<sup>19,25</sup>. Protein abundances are also less affected than mRNA levels by genetic  
48 variation<sup>29,30</sup>, including variations in gene copy numbers<sup>31–33</sup>. Consequently, protein  
49 expression profiling outperforms mRNA expression profiling with regard to gene function  
50 prediction<sup>19,20</sup>. Protein-based profiling not only allows for a more accurate measurement of  
51 gene activity, but can determine additional aspects of a cell's response to a perturbation,  
52 such as changes in protein localization and modification state. At the proteome level,  
53 expression profiling can therefore be extended to a more comprehensive protein covariation  
54 analysis.

55 Proof-of-principle studies by us and others have shown that protein covariation can  
56 be used to infer, for example, the composition of protein complexes and organelles<sup>34–42</sup>.  
57 However, these studies have focussed on relatively small sets of proteins or biological  
58 conditions, or used samples tailored to the analysis of specific cellular structures. In addition  
59 to the limited amount of data, coexpression analyses may be held back by the statistical  
60 tools used to pinpoint genes with similar activity. Coexpressed genes are commonly  
61 identified using Pearson's correlation, which is restricted to linear correlations and  
62 susceptible to outliers. Machine-learning may offer an increase in sensitivity and specificity.

63 Despite the success of functional genomics, many human proteins remain  
64 uncharacterized, especially small proteins that are difficult to study by biochemical methods.  
65 The emergence of big proteomics data and new computational approaches could provide an  
66 opportunity to look at these proteins from a different angle. We wondered if protein  
67 covariation would assign functions to previously uncharacterized proteins or novel roles to  
68 characterized ones. The resulting resource is available at [www.proteomeHD.net](http://www.proteomeHD.net) to generate  
69 hypotheses on the cellular functions of proteins of interest in a straightforward manner.

## 70 RESULTS

### 71 ProteomeHD is a data matrix for functional proteomics

72 To turn protein covariation analysis into a system-wide, generally applicable method, we  
73 created ProteomeHD. In contrast to previous drafts of the human proteome<sup>8,9,22,43,44</sup>,  
74 ProteomeHD does not catalogue the proteome of specific tissues or subcellular  
75 compartments. Instead, ProteomeHD catalogues the transitions between different proteome  
76 states, i.e. changes in protein abundance or localization resulting from cellular perturbations.  
77 HD, or high-definition, refers to two aspects of the dataset. First, all experiments are  
78 quantified using SILAC (stable isotope labelling by amino acids in cell culture)<sup>45</sup>. SILAC  
79 essentially eliminates sample processing artifacts and is especially accurate when

80 quantifying small fold-changes. This is crucial to detect subtle, system-wide effects of a  
81 perturbation on the protein network. Second, HD refers to the number of observations  
82 (pixels) available for each protein. As more perturbations are analysed, regulatory patterns  
83 become more refined and can be detected more accurately.

84 To assemble ProteomeHD we processed the raw data from 5,288 individual  
85 mass-spectrometry runs into one coherent data matrix, which covers 10,323 proteins (from  
86 9,987 genes) and 294 biological conditions (Supplementary Table 1). About 20% of the  
87 experiments were performed in our laboratory and the remaining data were collected from  
88 the Proteomics Identifications (PRIDE)<sup>46</sup> repository (Fig. 1a). The data cover a wide array of  
89 quantitative proteomics experiments, such as perturbations with drugs and growth factors,  
90 genetic perturbations, cell differentiation studies and comparisons of cancer cell lines  
91 (Supplementary Table 2). All experiments are comparative studies using SILAC<sup>45</sup>, i.e. they  
92 do not report absolute protein concentrations but highly accurate fold-changes in response to  
93 perturbation. About 60% of the included experiments analysed whole-cell samples. The  
94 remaining measurements were performed on samples that had been fractionated after  
95 perturbation, e.g. to enrich for chromatin-based or secreted proteins. This allows for the  
96 detection of low-abundance proteins that may not be detected in whole-cell lysates.

### 97 **ProteomeHD offers high protein coverage**

98 On average, the 10,323 human proteins in ProteomeHD were quantified on the basis of 28.4  
99 peptides and a sequence coverage of 49% (Supplementary Fig. 1). As expected from  
100 shotgun proteomics data, not every protein is quantified in every condition. The 294 input  
101 experiments quantify 3,928 proteins on average. Each protein is quantified, on average, in  
102 112 biological conditions (Supplementary Fig. 1). As a rule of thumb, coexpression studies  
103 discard transcripts detected in less than half of the samples. However, with 294 conditions  
104 ProteomeHD is considerably larger than the typical coexpression analysis. We therefore  
105 decided to use a lower arbitrary cut-off and include proteins for downstream analysis if they  
106 were quantified in about a third of the conditions. Specifically, we focus our co-regulation  
107 analysis on the 5,013 proteins that were quantified in at least 95 of the 294 perturbation  
108 experiments. On average, these 5,013 proteins were quantified in 190 conditions; 43% were  
109 quantified in more than 200 conditions (Supplementary Fig. 1).

### 110 **Machine-learning captures functional protein associations**

111 Proteins that are functioning together have similar patterns of up- and down regulation  
112 across the many conditions and samples in ProteomeHD. For example, the patterns of  
113 proteins belonging to two well-known biological processes, oxidative phosphorylation and  
114 rRNA processing, can be clearly distinguished, even though most expression changes are  
115 well below 2-fold (Fig. 1b). Therefore, we reasoned that it should be possible to reveal  
116 functional links between proteins on the basis of such regulatory patterns, and reveal the  
117 function of unknown proteins by associating them with well-characterized ones.

118 Traditionally, the extent of coexpression between two genes is determined by  
119 correlation analysis, for example using Pearson's correlation coefficient (PCC). Since PCC is  
120 very sensitive to outlier measurements, Spearman's rank correlation ( $\rho$ ) or Biweight  
121 midcorrelation (bicor) are sometimes used as more robust alternatives. We calculated these

122 three correlation coefficients for all 12,562,578 pairwise combinations of the 5,013 protein  
123 subset of ProteomeHD. To assess which metric works best for ProteomeHD we performed a  
124 precision-recall analysis, using known functional protein - protein associations from  
125 Reactome<sup>47</sup> as gold standard. This showed no major difference between the correlation  
126 measures, although Spearman's rho performs slightly better than the others (Fig. 1c).

127 We then tested a new type of coexpression measure based on unsupervised  
128 machine-learning. Specifically, we used the treeClust algorithm developed by Buttrey and  
129 Whitaker, which infers dissimilarities based on decision trees<sup>48,49</sup>. In short, treeClust runs  
130 data through a set of decision trees, which it creates without explicitly provided training data,  
131 and essentially counts how often two proteins end up in the same leaves. This results in  
132 pairwise protein - protein dissimilarities (not clusters of proteins). Importantly, we find that  
133 treeClust dissimilarities strongly outperform the three correlation metrics at predicting  
134 functional relationships between proteins in ProteomeHD (Fig. 1c).

135 Finally, we apply a topological overlap measure (TOM)<sup>50,51</sup> to the treeClust  
136 similarities, which further enhances performance by approximately 10% as judged by the  
137 area under the precision-recall curve (Fig. 1c). The TOM is typically used to improve the  
138 robustness of correlation networks by re-weighting connections between two nodes  
139 according to how many shared neighbors they have. The TOM-optimised treeClust results  
140 form our "co-regulation score". This score is continuous and reflects how similar two proteins  
141 behave across ProteomeHD, i.e. the higher the score the more strongly co-regulated two  
142 proteins are. However, for some questions a simplified categorical interpretation is more  
143 straightforward. In these cases we arbitrarily consider the top-scoring 0.5% percent of  
144 proteins pairs as "co-regulated". In this way, we identify 62,812 co-regulated protein pairs  
145 (Fig. 1d, Supplementary Table 3). For comparison, if the same data were analysed by  
146 Pearson's correlation, selecting the top 0.5% pairs would correspond to a cut-off of PCC >  
147 0.69, which is generally considered a strong correlation.

148 We then tested whether co-regulation indicates co-function. Indeed, we find that  
149 co-regulated protein pairs are heavily enriched for subunits of the same protein complex,  
150 enzymes catalysing consecutive metabolic reactions and proteins occupying the same  
151 subcellular compartments (Fig. 1e). The majority of proteins are co-regulated with at least  
152 one other protein, and about a third have more than five co-regulation partners (Fig. 1f). For  
153 99% of the tested proteins that had  $\geq 10$  co-regulated pairs, the group of their co-regulation  
154 partners is enriched in at least one Gene Ontology<sup>52</sup> biological process (Fig. 1g).

### 155 **Quantitative protein co-regulation is more informative than co-occurrence**

156 While decision trees are well-understood building blocks of many established  
157 machine-learning algorithms, treeClust itself is a relatively recent invention<sup>48</sup>. It was therefore  
158 unclear which type of information treeClust captures from a dataset. For example, treeClust  
159 scores could simply reflect whether or not two proteins are detected in the same set of  
160 samples, a measure that has been successfully exploited previously<sup>41</sup>. To test that we  
161 compared treeClust scores to the Jaccard index<sup>53</sup>, a dedicated measure of co-occurrence  
162 (Supplementary Fig. 2). In addition, we forced treeClust to learn dissimilarities solely based  
163 on co-occurrence by using a "binary" version of ProteomeHD, where all SILAC ratios were  
164 turned into ones and all missing values into zeroes. We find that the Jaccard index and

165 “binary” treeClust detect functionally related proteins equally well, but with much lower  
166 precision than standard treeClust. This suggests that protein co-regulation, i.e. coordinated  
167 changes in protein abundance, rather than co-detection is essential for treeClust  
168 performance.

169 Furthermore, it remained unclear what type of quantitative relationships treeClust can  
170 identify and why it outperforms correlation metrics for protein coexpression analysis. We  
171 addressed this in a separate study by systematically benchmarking treeClust using synthetic  
172 data<sup>54</sup> (available at: [www.biorxiv.org/content/10.1101/578971v1](http://www.biorxiv.org/content/10.1101/578971v1)). In short, we found that  
173 treeClust detects linear but not non-linear relationships. Unlike correlation metrics, it  
174 distinguishes between strong, tight-fitting relationships and weak trends. Finally, as may be  
175 expected from an algorithm based on decision trees, it is exceptionally robust against  
176 outliers. These properties of treeClust collectively explain its superior performance on  
177 ProteomeHD<sup>54</sup>. However, experiments with synthetic data also show that treeClust works  
178 best for large datasets with 50 samples or more, depending on additional parameters such  
179 as the frequency of missing values. Traditional correlation analysis may be better suited for  
180 smaller gene expression datasets<sup>54</sup>.

## 181 **A co-regulation map of the human proteome**

182 As a result of treeClust learning we know for each protein how strongly - or weakly - it is  
183 co-regulated with any other protein. In principle, these results could be displayed as a  
184 scale-free protein interaction network with edges indicating co-regulation (Supplementary  
185 Fig. 3). However, due to size and nature of our co-regulation data - 62,812 top-scoring links  
186 between 5,013 proteins - it appears impossible to avoid low-informative “hairball” graphs<sup>55</sup>.

187 We therefore chose to visualize the protein - protein co-regulation matrix using  
188 t-Distributed Stochastic Neighbor Embedding (t-SNE)<sup>56,57</sup>. This produces a two-dimensional  
189 proteome co-regulation map in which the distance between proteins indicates how similar  
190 they responded to the various perturbations in ProteomeHD (Fig. 1h, Supplementary Table  
191 4). Notably, t-SNE takes all pairwise co-regulation scores into account, rather than focussing  
192 on a small number of links above an arbitrary threshold. The t-SNE map shows that protein  
193 co-regulation is closely related to co-function. From a global perspective, the map reflects  
194 the subcellular organization of the cell (Fig. 1i). It broadly separates organelles and, for  
195 example, sets apart the nucleolus from the nucleus. A closer look into three sections of the  
196 map reveals that it captures more detailed functional relationships, too. For example, the five  
197 protein complexes of the respiratory chain are almost resolved (Fig. 1i, section 1). The  
198 section also contains the phosphate and ADP carriers that transport the substrates for ATP  
199 synthesis through the inner mitochondrial membrane, and ATP1F1 - a short-lived,  
200 post-transcriptionally controlled key driver of oxidative phosphorylation in mammals<sup>58</sup>.  
201 Similarly, cytoskeleton proteins such as actins and myosins are found next to their  
202 regulators, including Rho GTPases and the Arp2/3 complex (Fig. 1i, section 2). A third  
203 example section shows groups of proteins involved in RNA biology, from nucleolar rRNA  
204 processing to mRNA splicing and export (Fig. 1i, section 3). Notably, these annotations are  
205 only used to illustrate that the co-regulation map reflects functional similarity; the map itself is  
206 generated without any curated information, solely on the basis of protein abundance

207 changes in ProteomeHD. Therefore, the co-regulation map provides a data-driven overview  
208 of the proteome, connecting proteins into functionally related groups.

### 209 **Co-regulation complements existing functional genomics methods**

210 We next asked if protein co-regulation can predict associations that are not detected by other  
211 methods. For this we compare co-regulation to four alternative large-scale resources:  
212 IntAct<sup>59</sup>, BioGRID<sup>60</sup>, STRING<sup>61</sup> and BioPlex<sup>4</sup>. The first three are “meta-resources”, i.e. they  
213 compile curated sets of protein - protein interactions (PPIs) from the results of thousands of  
214 individual studies. Since meta-resources generally map interactions to gene loci rather than  
215 proteins, we disregard protein isoforms for this comparison and focus on co-regulated genes.

216 The co-regulation map covers fewer distinct genes than the other resources, but only  
217 STRING captures more interactions per average gene (Fig. 2a). Based on the 2,565 genes  
218 covered by both approaches, around 39% of the gene pairs identified as co-regulated had  
219 previously been linked in STRING (Fig. 2b). This suggests that co-regulation analysis  
220 confirms existing links, but also provides many additional ones. Conversely, only 7% of  
221 STRING PPIs are co-regulated, which may reflect the diverse molecular nature of  
222 associations covered by STRING. Notably, the overlap between the resources depends on  
223 the stringency setting: considering fewer, more stringent STRING interactions decreases the  
224 coverage of co-regulated genes and increases STRING PPIs identified as co-regulated (Fig.  
225 2b). An equivalent trend would be observed when modulating the co-regulation cut-off.  
226 STRING associations are based on multiple types of evidence, of which “mRNA  
227 coexpression” unsurprisingly shows the highest individual overlap with protein co-regulation  
228 results (Fig. 2c).

229 Next, we compared co-regulation specifically to physical PPIs catalogued by IntAct  
230 and BioGRID. We find that 11% of co-regulated gene pairs have a documented physical  
231 interaction between their proteins in BioGRID, and 3% are found in the smaller IntAct  
232 database (Fig. 2b). These physical PPIs were mainly derived from co-fractionation  
233 experiments, which tend to capture indirect interactions, rather than methods that detect  
234 direct interactions, such as two-hybrid screens (Fig. 2c).

235 Finally, we compared the co-regulation approach to an individual functional genomics  
236 project: BioPlex 2.0, the most comprehensive affinity purification–mass spectrometry  
237 (AP-MS) study reported to date<sup>4</sup>. BioPlex reports 4,935 physical interactions between the  
238 proteins used in our study, of which 19% are also co-regulated (Fig. 2d). An additional  
239 43,759 potential links between these proteins are identified uniquely by co-regulation. These  
240 are strongly enriched for functional protein associations found in STRING, compared to a  
241 random set of protein pairs (Fig. 2d). In conclusion, these comparisons suggest that protein  
242 co-regulation identifies protein - protein associations in a way that is reliable yet  
243 complementary to existing functional genomics methods. Note that proteins can interact  
244 physically or genetically or co-localize without being co-regulated, and vice versa. Therefore,  
245 protein co-regulation is complementary not just in terms of identifying new links, but also in  
246 providing additional, independent biological evidence for associations detected by other  
247 approaches.

## 248 **Uncharacterized proteins in ProteomeHD are rich in microproteins**

249 The co-regulation map contains 301 uncharacterized proteins, which we define as proteins  
250 with a UniProt<sup>62</sup> annotation score of 3 or less (Fig. 2e). Of these, 51% are co-regulated with  
251 at least one fully characterized protein, i.e. a protein with an annotation score of 4 or 5 (Fig.  
252 2f). On median, these uncharacterized proteins have 9 well-studied co-regulation partners,  
253 making it possible to predict their potential function in a “guilt by association” approach. We  
254 observe a similar connectivity for the cancer gene census<sup>63</sup>, i.e. genes that cause cancer  
255 when mutated, and for DisGeNET<sup>64</sup> genes, which are genes implicated in a broad range of  
256 human diseases (Fig. 2f). Therefore, protein co-regulation may also be helpful for functional  
257 analysis of human disease genes.

258 A common property of uncharacterized proteins is their small size. For example,  
259 proteins smaller than 15 kDa constitute 18% of the uncharacterized proteins in the human  
260 proteome, but only 5% of the characterized ones. Among the least well understood fraction  
261 of the proteome, i.e. proteins with an annotation score of 1, 40% are smaller than 15 kDa  
262 (Fig. 2g). This discrepancy is set to increase further, since hundreds or thousands such  
263 microproteins have so far been overlooked by genome annotation efforts<sup>65,66</sup>. Microproteins  
264 can regulate fundamental biological processes<sup>67</sup>, but their small size makes it difficult to  
265 identify interaction partners<sup>65,68</sup> or to target them in mutagenesis screens<sup>65</sup>. Microprotein  
266 sequences also tend to be less conserved than those of longer protein-coding genes<sup>69</sup>. We  
267 reasoned that our perturbation proteomics approach may help to reduce the annotation gap  
268 for small proteins. As it only requires proteins to be quantifiable in cell extracts we expect it  
269 to be less biased by protein size than methods involving extensive genetic or biochemical  
270 sample processing. Indeed, we find that 16% of the uncharacterized proteins in the  
271 co-regulation map are smaller than 15 kDa, which is close to the 18% in the proteome  
272 overall (Fig. 2h). However, it is a significant difference to BioPlex’s cutting-edge AP-MS data,  
273 in which microproteins drop to 6% ( $p < 2e-5$  in a one-tailed Fisher’s Exact test).

274 The fact that microproteins are not underrepresented in ProteomeHD does not  
275 automatically mean that their detection and characterisation is as robust as that of larger  
276 proteins. However, the average microprotein in the co-regulation map has been identified by  
277 12.2 peptides, many of which overlap and together result in an average sequence coverage  
278 of 76.4% (Supplementary Fig. 4a, d). While in a typical SILAC experiment proteins are  
279 considered to be quantifiable from upwards of two independent observations (SILAC ratio  
280 counts), microproteins in the co-regulation map are quantified with an average of 9 ratio  
281 counts per experiment, totalling a median of 671 ratio counts across ProteomeHD  
282 (Supplementary Fig. 4b, c). This indicates that microprotein quantitation in ProteomeHD is  
283 robust. Surprisingly, we find that microproteins have more co-regulation partners than larger  
284 proteins, and the same is true for their connectivity in STRING (Supplementary Fig. 4f).  
285 Within STRING, the majority of microprotein interactions are derived from curated  
286 annotations rather than high-throughput efforts such as RNA coexpression and text mining  
287 (Supplementary Fig. 4g). Note that, based on BioGRID, microproteins engage in fewer  
288 physical PPIs than larger proteins. This may be the result of an experimental bias  
289 (microproteins may dissociate more easily during purification and are more difficult to detect)  
290 or reflect a biological property (microproteins may have fewer physical interaction partners).

291 In either case, co-regulation offers itself as a powerful alternative approach to study  
292 microprotein functions in a systematic way.

### 293 **Functional annotation of proteins by co-regulation**

294 To facilitate the characterization of proteins through co-regulation we created the website  
295 [www.proteomeHD.net](http://www.proteomeHD.net). It allows users to search for a protein of interest, showing its position  
296 in the co-regulation map together with any co-regulation partners (Supplementary Fig. 5).  
297 The online map is interactive and zoomable, making it easy to explore the neighborhood of a  
298 query protein. The co-regulation score cut-off can be adjusted and statistical enrichment of  
299 Gene Ontology<sup>52</sup> terms among the co-regulated proteins is automatically calculated.

300 For example, protein co-regulation can be used to predict the potential function of  
301 uncharacterized microproteins such as the mitochondrial proteolipid MP68. MP68 is  
302 co-regulated with subunits of the ATP synthase complex, suggesting a function in ATP  
303 production (Fig. 1i, section 1). Despite being only 6.8 kDa small, its presence in the  
304 co-regulation map is documented by 8 distinct peptides that were observed a total of 398  
305 times across 142 experiments (Supplementary Fig. 4e). Intriguingly, MP68 co-purifies  
306 biochemically with the ATP synthase complex, but only in buffers containing specific  
307 phospholipids<sup>70,71</sup>, and knockdown of MP68 decreases ATP synthesis in HeLa cells<sup>72</sup>.

308 Virtually nothing is known about the 12 kDa microprotein TMEM256, although  
309 sequence analysis suggests it may be a membrane protein. Its position in the co-regulation  
310 map (Fig. 2i) and GO analysis of its co-regulation partners indicates that it likely localizes to  
311 the inner mitochondrial membrane (GO:0005743, Bonferroni adj.  $p < 5e-40$ ), where it may  
312 participate in oxidative phosphorylation (GO:0006119,  $p < 3e-35$ ).

313 Some proteins have no co-regulation partners above the default score cut-off, but  
314 can still be functionally annotated through the co-regulation map. The uncharacterized 224  
315 kDa protein HEATR5B, for example, is located in an area related to vesicle biology (Fig. 2i).  
316 Its immediate neighbours are five subunits of the HOPS complex, which mediates the fusion  
317 of late endosome to lysosomes. The position in the map shows that the HOPS complex is  
318 the closest fit to HEATR5B's regulation pattern, but they are not as similar as the top-scoring  
319 pairs in our overall analysis. If the co-regulation score cut-off is lowered, HOPS subunits and  
320 other endolysosomal proteins are eventually identified as co-regulated with HEATR5B, with  
321 concomitant enrichment of the related GO terms. This suggests that HEATR5B may not itself  
322 be a HOPS subunit, but could have a related vesicle-based function. Notably, a biochemical  
323 fractionation profiling approach also predicted HEATR5B to be a vesicle protein<sup>73</sup>.

324 Multifunctional proteins appear to fall into two categories in terms of co-regulation  
325 behavior. Prohibitin, for example, functions both as a mitochondrial scaffold protein and a  
326 nuclear transcription factor<sup>74</sup>. However, only the mitochondrial function is represented in the  
327 co-regulation map (Fig. 2j). This could indicate that its nuclear activity is not relevant in the  
328 biological conditions covered by ProteomeHD, or that only a small intracellular pool of  
329 prohibitin is nuclear, so that changes in its nuclear abundance are insignificant in comparison  
330 to the mitochondrial pool. In contrast, the helicase DDX3X shuttles between nucleus and  
331 cytoplasm, functioning both as nuclear mRNA processing factor and cytoplasmic regulator of  
332 translation<sup>75</sup>. In the co-regulation map, DDX3X sits between the areas related to these two  
333 activities and is significantly co-regulated both with proteins involved in nuclear RNA biology



334 and with translation factors (Fig. 2j). Therefore, DDX3X is a multifunctional protein whose  
335 separate activities result in a mixed regulatory pattern.

336 The protein co-regulation data presented here have been integrated into the recently  
337 released 11th version of STRING<sup>76</sup> (<https://string-db.org/>). In STRING's human protein -  
338 protein association network, links between proteins inferred from co-regulation in  
339 ProteomeHD are shown as network edges of the "coexpression" type (Supplementary Fig.  
340 6). Therefore, STRING is an alternative source for users wishing to explore protein  
341 co-regulation in conjunction with other types of association evidence.

### 342 **A new function for PEX11 $\beta$ in peroxisome-mitochondria interplay**

343 Some well-characterized proteins have unexpected co-regulation partners. For example,  
344 PEX11 $\beta$  is a key regulator of peroxisomal membrane dynamics and division<sup>77</sup>. However,  
345 PEX11 $\beta$ 's co-regulation partners are not peroxisomal proteins but subunits of the  
346 mitochondrial ATP synthase and other components of the electron transport chain (Fig. 1i,  
347 section 1). These proteins are located to the inner mitochondrial membrane, making a  
348 physical interaction with PEX11 $\beta$  unlikely. However, peroxisomes and mitochondria in  
349 mammals are intimately linked cooperating in fatty acid  $\beta$ -oxidation and ROS homeostasis<sup>78</sup>.  
350 How these organelles communicate or mediate metabolite flux has been elusive. Live cell  
351 imaging revealed that expression of PEX11 $\beta$ -EGFP in mammalian cells induced the  
352 formation of peroxisomal membrane protrusions, which interact with mitochondria (Fig. 3,  
353 Supplementary movies 1-3). Interactions of elongated peroxisomes with mitochondria were  
354 more frequent than those of spherical organelles, but both interactions were long-lasting  
355 (Fig. 3n,o). This indicates that peroxisome elongation can facilitate organelle interaction, but  
356 once organelles are tethered, the duration of contacts is similar between different  
357 morphological forms. Miro1 (RHOT1), a membrane adaptor for the microtubule-dependent  
358 motors kinesin and dynein<sup>79</sup>, is also co-regulated with PEX11 $\beta$  (Fig. 1i, section 1). We and  
359 others recently showed that Miro1 distributes to mitochondria and peroxisomes<sup>80,81</sup> indicating  
360 that it coordinates mitochondrial and peroxisomal dynamics with local energy turnover.  
361 Peroxisome-targeted Miro1 (Myc-Miro-PO) can be used as a tool to exert pulling forces at  
362 peroxisomal membranes, which results in the formation of membrane protrusions in certain  
363 cell types<sup>82</sup> (Supplementary Fig. 7). We show here that silencing of PEX11 $\beta$  inhibits  
364 membrane elongation by Myc-Miro-PO, confirming that PEX11 $\beta$  is required for the formation  
365 of peroxisomal membrane protrusions (Supplementary Fig. 7). These findings are in  
366 agreement with studies in plants, where *At*PEX11a has been reported to mediate the  
367 formation of peroxisomal membrane extensions in response to ROS<sup>83</sup>. In yeast,  
368 peroxisome-mitochondria contact sites are established by ScPex11 and ScMdm34, a  
369 component of the ERMES complex<sup>84</sup>. Additional tethering functions for the yeast mitofusin  
370 Fzo1 and ScPex34 in peroxisome-mitochondria contacts have recently been revealed<sup>85</sup>.  
371 Importantly, the study also demonstrated a physiological role for peroxisome-mitochondria  
372 contact sites in linking peroxisomal  $\beta$ -oxidation and mitochondrial ATP generation by the  
373 citric acid cycle<sup>85</sup>. We conclude that PEX11 $\beta$  and Miro1 contribute to peroxisome membrane  
374 protrusions, which present a new mechanism of interaction between peroxisomes and  
375 mitochondria in mammals. They likely function in the metabolic cooperation and crosstalk  
376 between both organelles, and may facilitate transfer of metabolites such as acetyl-CoA

377 and/or ROS homeostasis during mitochondrial ATP production. These findings now enable  
378 future studies on the precise functions of peroxisome membrane protrusions in mammalian  
379 cells and the role of PEX11 $\beta$ .

### 380 **Proteomics enables higher accuracy but lower coverage than transcriptomics**

381 To compare the impact of mRNA and protein abundances on expression profiling we first  
382 focussed on 59 SILAC ratios in ProteomeHD that measured abundance changes across a  
383 panel of lymphoblastoid cell lines<sup>30</sup>. For these samples, corresponding mRNA abundance  
384 changes have been determined using RNA-sequencing<sup>86</sup>. Repeating treeClust learning on  
385 the basis of these data, we observed that protein coexpression predicts functional  
386 associations with far higher precision than mRNA coexpression (Fig. 4a). Similar results  
387 have recently been reported for a panel of human cancer samples<sup>19</sup>.

388 Such analyses show that in a direct gene-by-gene, sample-by-sample comparison,  
389 protein expression levels are better indicators for gene function than mRNA expression.  
390 However, the amount of transcriptomics data published to date vastly exceeds that of  
391 proteomics studies. For example, the NCBI GEO repository currently holds mRNA  
392 expression profiling data from more than one million human samples<sup>87</sup>. This raises the  
393 possibility that the sheer quantity of available transcriptomics data could overcome their  
394 reduced reflection of functional links and, in combined form, perform better than  
395 protein-based measurements. To test this we compared the ProteomeHD co-regulation  
396 score with Pearson correlation coefficients obtained by STRING, which leverages the vast  
397 amount of mRNA expression experiments deposited in GEO<sup>61,88</sup>. Remarkably,  
398 precision-recall analysis shows that the protein co-regulation score still outperforms mRNA  
399 coexpression, despite being based on only 294 SILAC ratios (Fig. 4b). Much of this  
400 improvement is due to the robustness of treeClust machine-learning, as Pearson's  
401 correlation coefficients derived from the same ProteomeHD data work only moderately better  
402 than mRNA correlation (Fig. 4b). While only gene pairs with both mRNA and protein  
403 expression measurements were considered for the precision-recall analysis, the  
404 transcriptomics and proteomics datasets individually covered 17,436 and 4,976 genes,  
405 respectively (Fig. 4b). Therefore, mRNA profiling outperforms protein profiling in terms of  
406 gene coverage. In addition, transcriptomics remains the only expression profiling approach  
407 suitable for non-coding RNAs.

## 408 **DISCUSSION**

409 ProteomeHD in conjunction with machine learning provides an entry point for "big-data"-type  
410 protein co-regulation analysis into the functional genomics methods repertoire. It is possible  
411 that accuracy and coverage could be increased further by adding additional proteomics data.  
412 To test this we randomly removed 5%, 10% or 15% of the data points in ProteomeHD. This  
413 decreases performance reproducibly and proportionally to the amount of removed data  
414 (Supplementary Fig. 8), suggesting that ProteomeHD has not reached saturation and  
415 expanding it will further enhance its performance. One possibility would be to incorporate  
416 other types of proteomics experiments, such as affinity-purifications or indeed the entire

417 PRIDE<sup>46</sup> repository. The latter approach is for instance taken by the Tabloid Proteome, which  
418 infers protein associations based on detecting them in the same subset of many different  
419 proteomics experiments<sup>41</sup>. However, there is a benefit of restricting ProteomeHD to  
420 perturbation experiments. It supports a biological interpretation of protein associations  
421 derived from it: two co-regulated proteins are part of the same cellular response to changing  
422 biological conditions, even though the precise molecular nature of the connection remains  
423 unknown. In this way, protein co-regulation analysis is analogous to genetic interaction  
424 screening. This also sets protein co-regulation apart from indiscriminate protein covariation  
425 or co-occurrence analyses, which find protein links in a mix of proteomics data and therefore  
426 give no insight into the possible biological connection.

427 A key difference between our approach and previous gene coexpression studies is  
428 our application of two machine-learning algorithms, treeClust<sup>48</sup> and t-SNE<sup>56,57</sup>. Inferring  
429 protein associations through treeClust learning is both more robust and sensitive than a  
430 traditional correlation-based approach, providing a leap in the accuracy with which  
431 functionally relevant interactions can be identified from the same dataset. For example, a  
432 recent study reported a protein co-regulation network across 41 cancer cell lines and  
433 subsequently identified dysregulated protein associations that predict drug sensitivities of  
434 these cell lines<sup>20</sup>. Applying Spearman's correlation to high-quality, TMT-based proteomics  
435 data allowed Lapek *et al*<sup>20</sup> to detect protein-protein associations with an accuracy that was  
436 tenfold higher than that based on matching mRNA coexpression data. When applying  
437 treeClust to these data, strikingly, we can further improve this performance (Supplementary  
438 Fig. 9a). This suggests that treeClust may be helpful for the detection of "dysregulation  
439 biomarkers" in the future. The second machine-learning tool we apply here, t-SNE, visualizes  
440 treeClust-learned protein associations as a 2D map. Correlation networks are typically built  
441 from a limited number of the strongest pairwise interactions, whereas t-SNE takes into  
442 account the similarity - or dissimilarity - between all possible pairwise protein combinations. It  
443 creates the map that best reflects both direct and indirect relationships between all proteins.  
444 In this way, also proteins that are not directly linked to the core network can be placed into a  
445 functional context. For example, a t-SNE co-regulation map obtained for Lapek *et al*'s cancer  
446 proteomics dataset contains the complete set of ~6,800 proteins, rather than the 3,024  
447 proteins that are directly correlated with another protein (Supplementary Fig. 9b). Moreover,  
448 protein-protein associations visualized by t-SNE can be explored in a hierarchical manner,  
449 with larger distances indicating weaker co-regulation. This may be useful for studying  
450 connections between related protein complexes (Fig. 1i) or to reveal broad functional clues  
451 for uncharacterized proteins for which no detailed predictions are available, such as the  
452 HEATR5B protein assigned to the vesicle area of the co-regulation map (Fig. 2i). Our web  
453 application at [www.proteomeHD.net](http://www.proteomeHD.net) is designed to support researchers in exploring  
454 co-regulation data at multiple scales, to validate existing hypotheses or create new ones.

455 Protein coexpression analysis identifies functional connections between proteins with  
456 an accuracy and sensitivity that is substantially higher than traditional mRNA coexpression  
457 analysis. This may be particularly important for constitutively active genes, which constitute  
458 about half of human genes<sup>44</sup> and are primarily controlled at the protein level<sup>89,90</sup>. With an ever  
459 increasing amount of protein expression data making their way into the public domain, and  
460 the simplicity of exploiting the analysis results by the scientific community, protein

461 coexpression analysis has a large potential for gene function annotation. Only 300  
462 quantitative proteomics measurements sufficed in conjunction with machine learning to  
463 establish functional connections between many human genes, which may be of considerable  
464 interest for proteome annotation in less studied or difficult to study organisms.

#### 465 **ACKNOWLEDGEMENTS**

466 We are grateful to Damian Szklarczyk for providing the mRNA Pearson correlation data used  
467 by STRING and the STRING team for testing our coregulation data and adding it as novel  
468 evidence type to STRING 11. We also thank Karen Wills, Kyosuke Nakamura, Constance  
469 Alabert and Anja Groth for contributing chromatin enrichment experiments, and Afsoon S.  
470 Azadi for support with live-cell-imaging. This work was supported by the Wellcome Trust  
471 through a Senior Research Fellowship to J.R. (grant number 103139) and by the  
472 Biotechnology and Biological Sciences Research Council (BB/N01541X/1, BB/R016844/1; to  
473 M.S.) and H2020-MSCA-ITN-2018 812968 PERICO (to M.S.). The Wellcome Centre for Cell  
474 Biology is supported by core funding from the Wellcome Trust (grant number 203149).

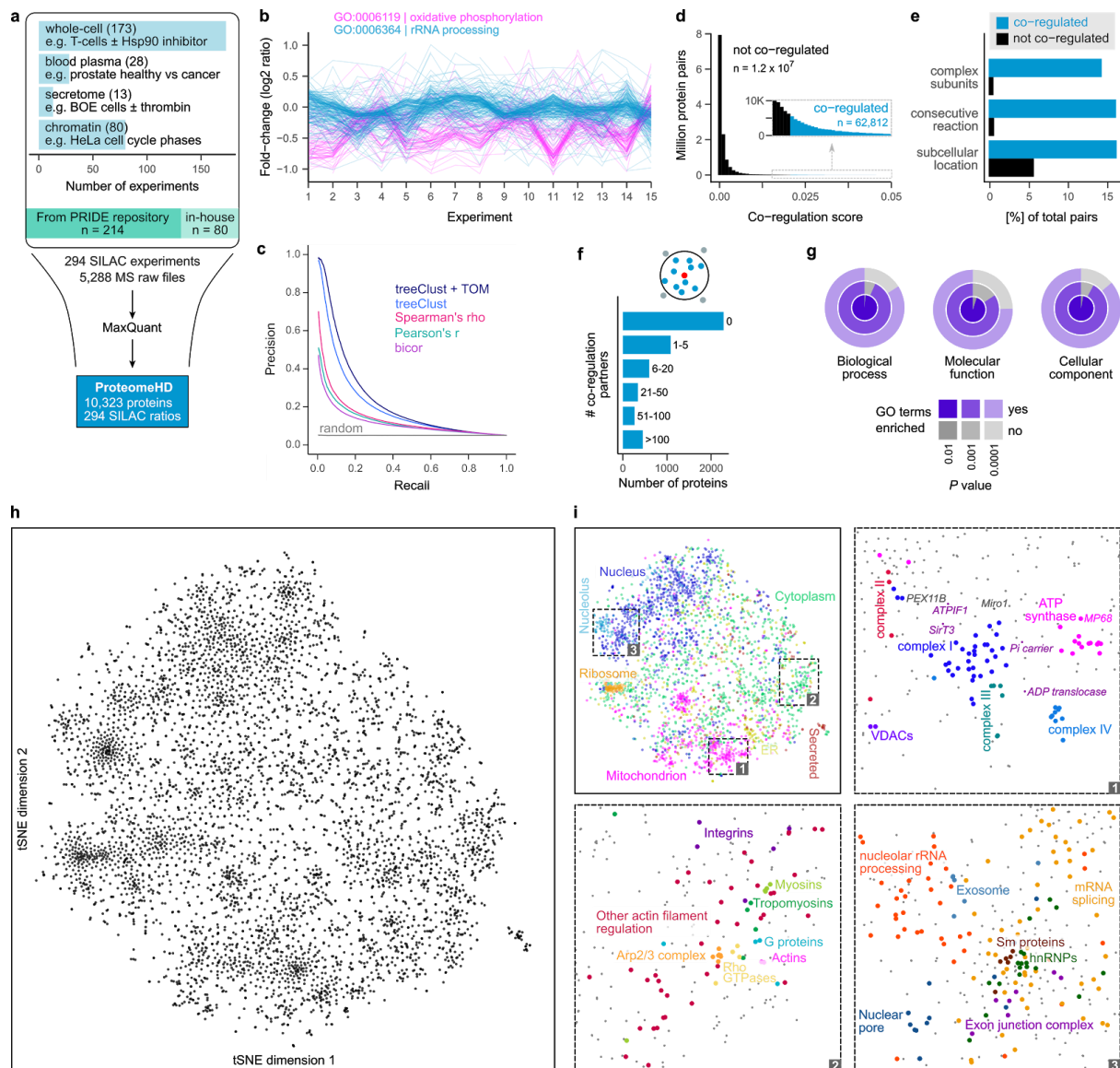
#### 475 **AUTHOR CONTRIBUTIONS**

476 G. K. and J. R. conceived the project. G. K. and P.G. conducted the data analysis. P. G.  
477 created the web application. T. A. S., J. B. P. and M. S. conducted the Pex11 $\beta$  analysis. All  
478 authors contributed to writing the manuscript.

#### 479 **COMPETING FINANCIAL INTERESTS**

480 The authors declare no competing financial interests.

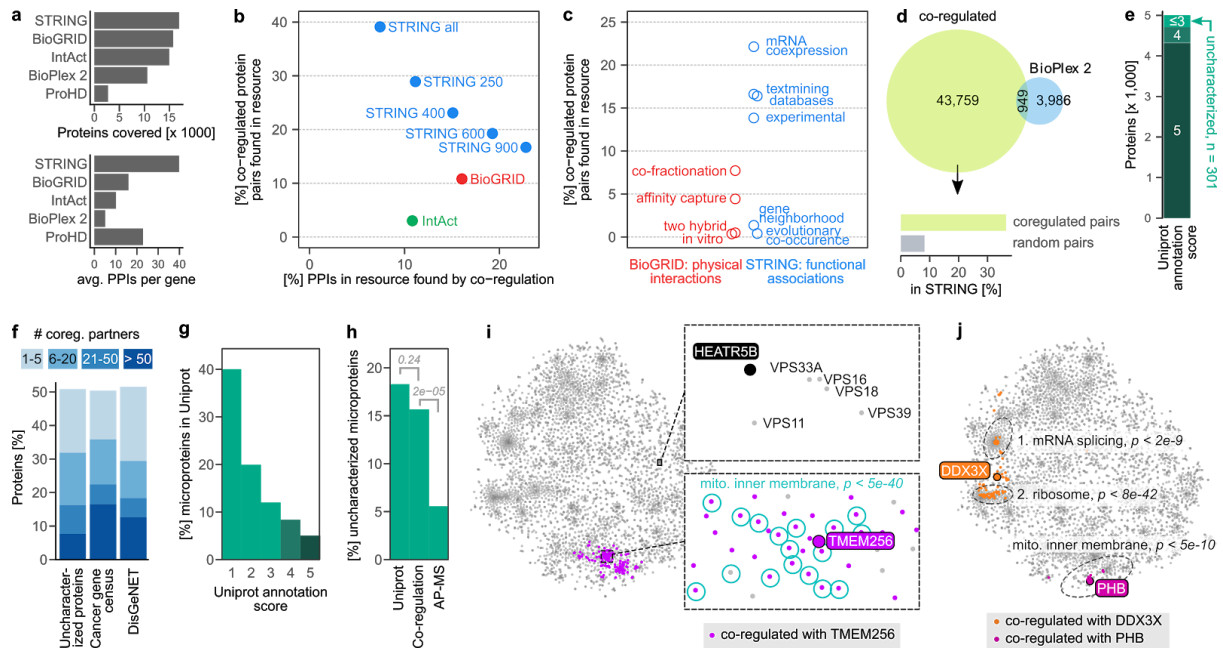
481 **MAIN FIGURES**



482 **Figure 1. The co-regulation map shows functional associations between human**  
483 **proteins.**

484 (a) Assembly of ProteomeHD, which quantifies the protein response to 294 perturbations  
485 using SILAC<sup>45</sup>. Most measurements document protein abundance changes in whole-cell  
486 samples, but in some cases subcellular fractions were enriched to detect low-abundance  
487 proteins. Data were collected from PRIDE<sup>46</sup> and produced in-house. (b) A random set of  
488 experiments from ProteomeHD, showing that groups of proteins with related functions, e.g.  
489 Gene Ontology<sup>52</sup> (GO) biological processes, display similar expression changes. Note that  
490 the fold-changes are often very small. (c) Precision - recall analysis showing that the  
491 treeClust<sup>48,49</sup> algorithm outperforms three correlation-based coexpression measures.  
492 Applying the topological overlap measure (TOM) improves performance further. Annotations  
493 in Reactome<sup>47</sup> were used as gold standard. (d) Co-regulation scores for all protein pairs are  
494 obtained by combining treeClust with TOM. The score distribution is highly skewed. Where

495 an arbitrary threshold is required, the highest-scoring 0.5% of pairs (N = 62,812) are  
496 considered “co-regulated”. **(e)** Co-regulated protein pairs are strongly enriched for subunits  
497 of the same protein complex, enzymes catalysing consecutive metabolic reactions and  
498 proteins with identical subcellular localization. **(f)** Most proteins are co-regulated with no or  
499 few other proteins, but many have more than 5 co-regulated partners. **(g)** Considering  
500 proteins that are co-regulated with  $\geq 10$  proteins, these groups of co-regulated proteins are  
501 almost always enriched in one or more GO terms. **(h)** The global co-regulation map of  
502 ProteomeHD created using t-Distributed Stochastic Neighbor Embedding (t-SNE)<sup>56,57</sup>.  
503 Distances between proteins indicate how similar their expression patterns are. See  
504 [www.proteomeHD.net](http://www.proteomeHD.net) for an interactive version of the map. **(i)** The co-regulation map  
505 broadly corresponds to subcellular compartments, and more detailed functional associations  
506 can be observed at higher resolution, as exemplified in subpanels 1-3.



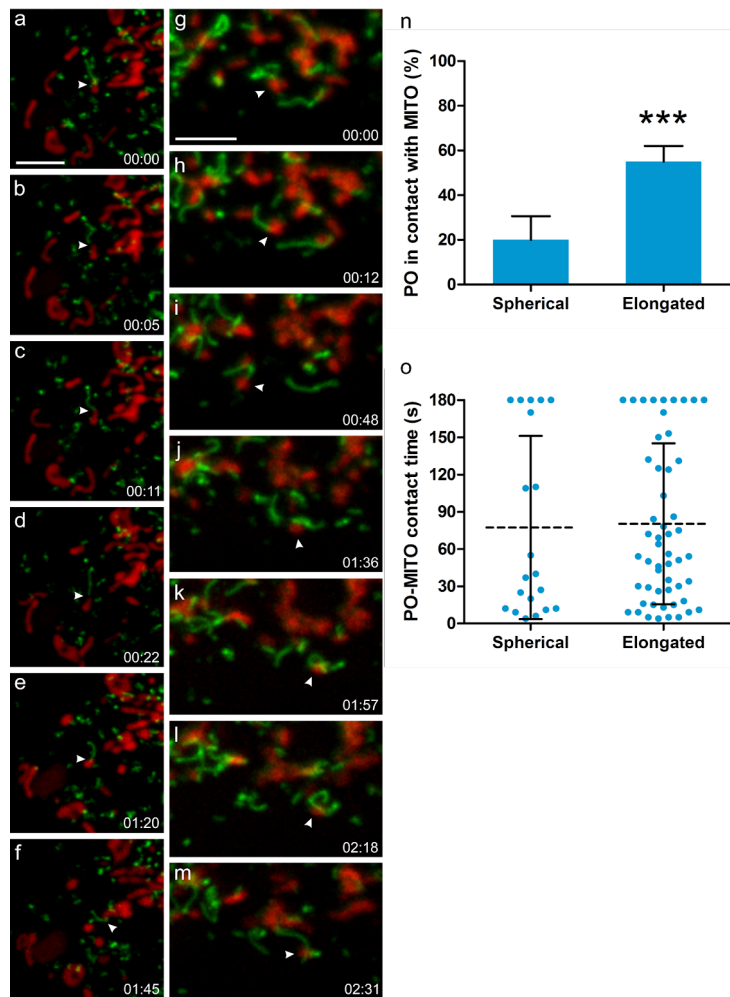
507 **Figure 2. Protein co-regulation complements existing methods and predicts functions**  
 508 **of unknown proteins.**

509 (a) Coverage of protein - protein interactions (PPIs) in comparison to other resources. Top  
 510 barchart shows the number of genes covered, i.e. having at least one PPI above cut-off.  
 511 STRING cut-off used: medium (400). Bottom chart shows the average number of PPIs of  
 512 covered genes. The co-regulation map (ProHD) covers fewer genes than STRING, BioGRID,  
 513 IntAct and BioPlex 2, but covers many associations between those genes. (b) Overlap  
 514 between PPIs discovered by protein co-regulation and PPIs already present in large-scale  
 515 annotation resources that cover both physical (BioGrid<sup>60</sup> and IntAct<sup>59</sup>) and functional  
 516 (STRING<sup>61</sup>) associations. Multiple association score cut-offs were considered for STRING.  
 517 These three resources integrate data from many small and large-scale studies. (c) Coverage  
 518 of co-regulated protein pairs in BioGRID and STRING broken down by the type of functional  
 519 genomics evidence available in each resource. (d) Number of co-regulation links compared  
 520 to PPIs found for the same set of genes by BioPlex 2.0<sup>4</sup>, one of the largest PPI datasets  
 521 reported to date by a single study. Associations unique to co-regulation are strongly enriched  
 522 for links in STRING, compared to random gene pairs. (e) Out of the 5,013 proteins in the  
 523 co-regulation map, 301 have a UniProt annotation score  $\leq 3$  and are thus defined as  
 524 uncharacterized. (f) Connectivity of either uncharacterized proteins or proteins encoded by  
 525 disease genes to well-characterized proteins (annotation score  $\geq 4$ ). 51% of uncharacterized  
 526 proteins have at least one co-regulation partner, 32% have more than five. (g) Barchart  
 527 showing the percentage of all 20,408 human UniProt (SwissProt) proteins that are  
 528 microproteins, i.e. have a molecular weight  $< 15$  kDa. Note that microproteins are heavily  
 529 enriched among less well-characterized proteins. (h) 18% of uncharacterized proteins in  
 530 UniProt are microproteins, compared to 16% of the uncharacterized proteins in the  
 531 co-regulation map and 6% in state-of-the-art AP-MS experiments, represented by BioPlex.  
 532 *P*-values are from one-sided Fisher's Exact test. (i) The uncharacterized microprotein

533 TMEM256 has many co-regulation partners, which are enriched for GO term “mitochondrial  
534 inner membrane” among others. Bonferroni-adjusted  $P$ -value is from a hypergeometric test.  
535 The uncharacterized HEATR5B protein has no co-regulation partners above the default  
536 threshold, but its position in the map nevertheless indicates a potential function. (j) For  
537 multifunctional proteins, co-regulation can reveal a mix of their functions (DDX3X), or their  
538 main function only (prohibitin, PHB). Three representative GO terms are shown.

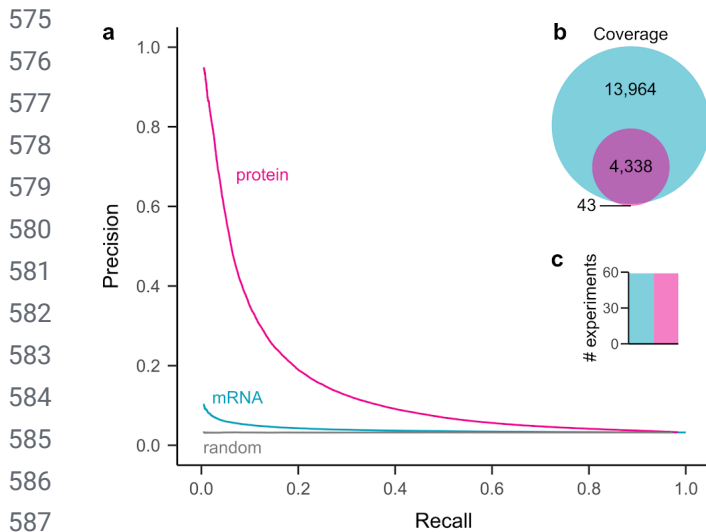


539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574



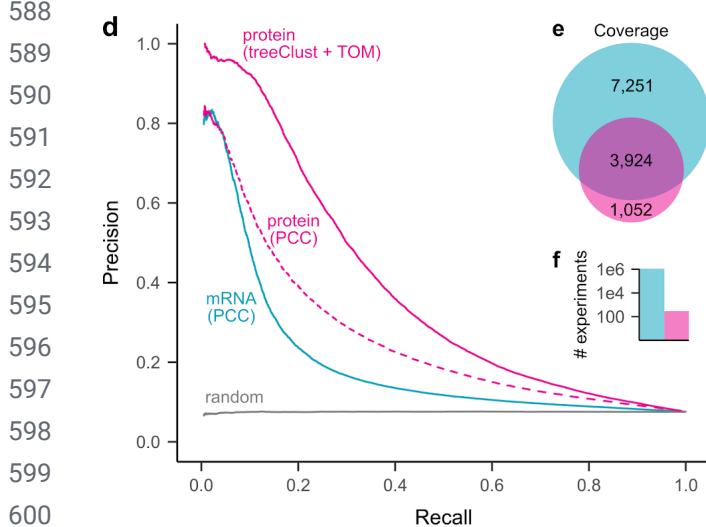
**Figure 3. PEX11 $\beta$  mediates the formation of peroxisomal membrane protrusions which interact with mitochondria in mammalian cells.**

(a-m) COS-7 cells were transfected with PEX11 $\beta$ -EGFP, mitochondria were stained with Mitotracker (red) and cells observed live using a spinning disc microscope. PEX11 $\beta$ , a membrane shaping protein, induces the formation of tubular membrane protrusions from globular peroxisomes. We show here that those membrane protrusions can interact with mitochondria. (a-f) shows a peroxisome which interacts with a mitochondrion via its membrane protrusion (arrowhead), and follows it, occasionally detaching and re-establishing contact before interacting with another mitochondrion (see Supplementary Movie 1). (g-m) shows a mitochondrion (arrowhead) which interacts with a peroxisome via a peroxisomal membrane protrusion. It then detaches and moves away to interact with another peroxisome, which wraps its protrusion around it, before interacting with another mitochondrion (see Supplementary Movie 2). (n) Quantification of interactions between spherical or elongated peroxisomes (PO) with mitochondria (MITO). The average result of 3 independent experiments is shown, error bars indicate standard deviation. (o) Quantification of contact time. Note that elongated PO interact more frequently with MITO than spherical PO, but for similar time periods. PO-MITO interactions are generally long-lasting (see Supplementary Movie 3) (n=200 peroxisomes from 5 different cells). Dotted line indicates the mean, error bars indicate standard deviation. \*\*\*  $P < 0.001$  from a two-tailed unpaired  $t$  test; Time (min:sec). Scale bars, 5  $\mu$ m.



**Figure 4. Protein co-regulation enables higher precision from less data, but has lower coverage than classic mRNA coexpression.**

(a) Precision-recall analysis of treeClust machine-learning on a subset of ProteomeHD, that is 59 samples for which matching RNA-seq data were available from a separate study<sup>86</sup>. Reactome pathways were used as gold standard for true functional associations (proteins found in same pathway) and false associations (never found in same pathway). Only annotated genes covered by both datasets were considered for PR analysis (n = 2,901). (b) Venn diagram showing number of genes covered by each analysis. (c) Bar chart showing number of experiments the curves are based on. (d) Similar precision-recall analysis of treeClust machine-learning on the full ProteomeHD database, in comparison to Pearson correlation obtained by STRING<sup>61</sup> on the basis of one million human mRNA profiling samples deposited in the NCBI Gene Expression Omnibus<sup>87</sup> ("mRNA / PCC"). Protein co-regulation outperforms mRNA



602 correlation despite being based on orders-of-magnitude less data. This is partially due to the use of machine-learning, as predicting associations from ProteomeHD using PCC decreases performance markedly ("protein / PCC"). Only annotated genes covered by both datasets were considered for the PR analysis (n = 2,743). (e, f) same as (b, c).

606 **SUPPLEMENTARY MOVIE LEGENDS**

607 **Supplementary Movie 1. Interaction of peroxisomal membrane protrusions with**  
608 **mitochondria in COS-7 cells. See Fig. 4a-f.**

609 COS-7 cells were transfected with PEX11 $\beta$ -EGFP, mitochondria were stained with  
610 Mitotracker (red), and analysed by live-cell imaging using an IX81 microscope (Olympus)  
611 equipped with a CSUX1 spinning disk head (Yokogawa). A peroxisome interacts with a  
612 mitochondrion via its membrane protrusion, and follows it, occasionally detaching and  
613 re-establishing contact. 200 stacks of 9 planes (0.5  $\mu$ m thickness, 100 ms exposure) were  
614 taken in a continuous stream. 118 frames, 14 $\times$  speed. Scale bar, 5  $\mu$ m.

615 **Supplementary Movie 2. Interaction of peroxisomal membrane protrusions with**  
616 **mitochondria in COS-7 cells. See Fig. 4g-m and legend Movie 1.**

617 Note a peroxisome at the bottom, which interacts with a mitochondrion via its membrane  
618 protrusion and then wraps around it, possibly to increase the membrane contact area. 200  
619 stacks of 9 planes (0.5  $\mu$ m thickness, 100 ms exposure) were taken in a continuous stream.  
620 200 frames, 14 $\times$  speed. Scale bar, 5  $\mu$ m.

621 **Supplementary Movie 3. Interaction of peroxisomal membrane protrusions with**  
622 **mitochondria in COS-7 cells. See legend Movie 1.**

623 A mitochondrion, which moves to the left, is dragging a peroxisome with a membrane  
624 protrusion with it, indicating that the organelles are tightly tethered to each other. 200 stacks  
625 of 9 planes (0.5  $\mu$ m thickness, 100 ms exposure) were taken in a continuous stream. 100  
626 frames, 14 $\times$  speed. Scale bar, 5  $\mu$ m.

## 627 ONLINE METHODS

### 628 General data analysis and code availability

629 Data analysis was performed in R<sup>91</sup>. R scripts and input files required to reproduce the  
630 results of this manuscript are available in the following GitHub repository:  
631 <https://github.com/Rappsilber-Laboratory/ProteomeHD>. The R package `data.table`<sup>92</sup> was  
632 used for fast data processing. Figures were prepared using `ggplot2`<sup>93</sup>, `gridExtra`<sup>94</sup>, `cowplot`<sup>95</sup>  
633 and `viridis`<sup>96</sup>.

### 634 Data selection for ProteomeHD

635 MS raw data were produced in-house or downloaded from the PRIDE repository<sup>46</sup>. Only  
636 experiments fulfilling the following inclusion criteria were considered:

637 (1) Comparative proteomics experiments, i.e. relative protein quantitations of two or  
638 more biological states. For example, cells treated with an inhibitor vs. mock control. (2)  
639 Biological - not biochemical - comparisons, i.e. fold-changes must have been brought about  
640 *in vivo*, not by differential biochemical purification. For example, SILAC-labelled cells were  
641 treated with inhibitor or mock control, harvested and combined, and chromatin was enriched  
642 on the combined sample. In such cases any observed fold-change reflects the response to  
643 the inhibitor in the living cell, for example a protein re-localising from cytoplasm onto  
644 chromatin. We did not consider experiments that compared, for example, a whole-cell lysate  
645 with a chromatin-enriched fraction, as this would measure the impact of the biochemical  
646 enrichment rather than a biological event. (3) Quantitation by “stable isotope labeling by  
647 amino acids in cell culture” (SILAC)<sup>45</sup>. (4) Samples of human origin.

648 In addition to these conceptual considerations, the following restrictions were  
649 imposed by the data processing pipeline: (5) The SILAC mass shift introduced by heavy  
650 arginine must be distinct from heavy lysine. (6) Raw data acquired on an Orbitrap mass  
651 spectrometer. (7) Samples alkylated with iodoacetamide, resulting in carbamidomethylation  
652 of cysteines.

653 In total, we considered 294 experiments (SILAC ratios) from 31 projects. A full list of  
654 these is provided in Supplementary Table 2, which also includes the PRIDE identifiers of all  
655 previously published datasets.

### 656 In-house data collection

657 80 experiments were performed in-house and analyzed chromatin-enriched samples. Of  
658 these, 65 measured the effect of growth factors, radiation and other perturbations on  
659 interphase chromatin, which was prepared using Chromatin Enrichment for Proteomics  
660 (ChEP)<sup>97</sup>. About half of these experiments had previously been published<sup>36</sup>. Another 15  
661 experiments documented perturbations specifically on freshly replicated chromatin, which  
662 was prepared using Nascent Chromatin Capture (NCC)<sup>98</sup>. All mass spectrometry raw files  
663 generated in-house have been deposited to the ProteomeXchange Consortium  
664 (<http://proteomecentral.proteomexchange.org>) via the PRIDE partner repository<sup>46</sup> with the  
665 dataset identifier PXD008888 (this repository will be made public upon acceptance of the  
666 manuscript).

## 667 **MS raw data processing**

668 The 5,288 MS raw files were processed using MaxQuant 1.5.2.8<sup>99</sup> on a Dell PowerEdge  
669 R920 server. The following default MaxQuant search parameters were used: MS1 tolerance  
670 for the first Andromeda search: 20 ppm, MS1 tolerance for the main Andromeda search: 4.5  
671 ppm, FTMS MS2 match tolerance: 20 ppm, ITMS MS2 match tolerance: 0.5 Da, Variable  
672 modifications: acetylation of protein N-termini, oxidation of methionine, Fixed modifications:  
673 carbamidomethylation of cysteine, Decoy mode set to reverse, Minimum peptide length: 7  
674 and Max missed cleavages set to 2. The following non-default settings were used: In  
675 group-specific parameters, match type was set to “No matching”. In global parameters,  
676 “Re-quantify” was enabled, minimum ratio count was set to 1 and “Discard unmodified  
677 counterpart peptide” was disabled. Also in global parameters, writing of large tables was  
678 disabled. SILAC labels were set as group-specific parameters as indicated in Supplementary  
679 Table 2. Canonical and isoform protein sequences were downloaded from UniProt<sup>62</sup> on 28th  
680 May 2015, considering only reviewed SwissProt entries that were part of the human  
681 proteome. Unprocessed MaxQuant result tables, including peptide evidence data, have been  
682 deposited into the PRIDE repository PXD008888.

683 Protein fold-changes were extracted from the MaxQuant proteinGroups file returned  
684 by MaxQuant. Non-normalized SILAC ratios were considered for downstream analysis, log2  
685 transformed and median-normalised. From triple labelling experiments, the heavy/light and  
686 medium/light ratios - but not the heavy/medium ratios - were considered. Proteins detected  
687 in less than 4 experiments were discarded, as were proteins labeled as contaminants,  
688 reverse hits and those only identified by a modification site. The resulting data matrix,  
689 ProteomeHD, can be downloaded as Supplementary Table 1.

## 690 **Calculation of treeClust dissimilarities**

691 It is common in gene coexpression studies to remove genes that were detected in less than  
692 half of the samples from the analysis. However, given the unusually large size of  
693 ProteomeHD we chose a different arbitrary cut-off, excluding proteins that were detected in  
694 less than 95 (about a third) of the 294 experiments. For the remaining 5,013 proteins in  
695 ProteomeHD we used the treeClust<sup>48</sup> R package to calculate all 12,562,578 pairwise  
696 dissimilarities. Note that treeClust was designed not only to measure inter-point  
697 dissimilarities but also to perform clustering<sup>48,49</sup>. However, in this study we use it only to  
698 calculate dissimilarities, via the treeClust.dist function. The dissimilarity specifier was set to  
699 d.num = 2, so that dissimilarities are weighted according to tree quality. We optimised two  
700 hyperparameters of treeClust and rpart, which is the routine treeClust uses to create  
701 decision trees. These were treeClust’s serule argument, which defines to extent to which  
702 trees are pruned, and rpart’s complexity (cp) parameter, which describes the improved fit  
703 required to attempt a split. A grid search was performed against the Reactome gold standard  
704 (see below) and the area under precision - recall curves was used to identify optimal  
705 parameter settings. They were determined to be serule = 1.8 and cp = 0.105, providing  
706 approximately a 10% performance improvement over treeClust’s default settings.

## 707 **Protein co-regulation scores**

708 To calculate the final pairwise co-regulation scores, treeClust dissimilarities were  
709 transformed further. First, they were turned into similarities, i.e.  $1 - \text{treeClust dissimilarity}$ .  
710 Using the WGCNA<sup>100,101</sup> R package, we then performed a sigmoid transformation of these  
711 treeClust similarities, creating an adjacency matrix. The settings of parameters  $\mu$  and  $\alpha$   
712 for this transformation were optimised in a grid search against the Reactome gold standard,  
713 using the area under precision - recall curves as readout. In a third step, the adjacency  
714 matrix was transformed into a topological overlap matrix using WGCNA's TOMsimilarity  
715 function, with the TOMDenom parameter set to "mean". These TOM similarities are the  
716 co-regulation scores used throughout our analysis. Co-regulation scores for all of our  
717 12,562,578 protein pairs can be downloaded from the PRIDE repository PXD008888.

718 While the co-regulation score is continuous, some analyses benefitted from a  
719 simplified categorical approach. For these cases we arbitrarily defined the highest-scoring  
720 0.5% of protein pairs as "co-regulated pairs" and the remaining 99.5% of pairs as "not  
721 co-regulated pairs". A list of all 62,812 co-regulated protein pairs is available as  
722 Supplementary Table 3.

## 723 **Reactome gold standard**

724 A gold standard set of reference proteins was defined using Reactome<sup>47</sup>. Bona fide  
725 functionally associated protein pairs (true positives) were defined as protein pairs found in  
726 the same "detailed" Reactome pathway. This was inferred from the file UniProt2Reactome.txt  
727 (available at <https://reactome.org/download-data>), where each protein is annotated to the  
728 lowest level subset of Reactome pathways. To make sure that only closely related protein  
729 pairs were assigned the "true positive" label, we excluded two pathways that were composed  
730 of > 200 proteins. We defined protein pairs that are not functionally associated (false  
731 positives) as proteins that are never in the same Reactome pathway, at any annotation level.  
732 This was inferred from UniProt2Reactome\_All\_Levels.txt (also available at  
733 <https://reactome.org/download-data>), a file that maps proteins to all levels of the Reactome  
734 pathway hierarchy. A copy of this gold standard is available in the Github repository noted  
735 above.

## 736 **Comparison of treeClust and correlation metrics**

737 Pearson's correlation coefficients (PCC) and Spearman's rank correlation coefficients ( $\rho$ )  
738 were obtained using the cor function in R, for the same protein pairs covered by the  
739 treeClust analysis. Biweight mid-correlation coefficients (bicor) were calculated with default  
740 settings using the R package WGCNA<sup>101,102</sup>. Changing the maxPOutliers parameter of the  
741 bicor function did not improve performance. Precision - recall (PR) analysis was performed  
742 with the ROCR package<sup>103</sup> using true and false positive pairs compiled from annotation in  
743 Reactome (see paragraph Reactome gold standard). The random classifier was created by  
744 scrambling co-regulation scores.

## 745 **t-SNE visualization**

746 To visualize ProteomeHD as a 2D co-regulation map, co-regulation scores were subjected to  
747 t-Distributed Stochastic Neighbor Embedding (t-SNE)<sup>56,57</sup> using the Rtsne<sup>104</sup> package for R.

748 The theta parameter was set to zero to calculate the exact embedding. The perplexity  
749 parameter was set to 50, up from the default of 30, to account for the large size of the  
750 co-regulation dataset. 1,500 iterations were performed. However, visual comparison of the  
751 t-SNE maps showed that these parameter adaptations provided only a marginal  
752 improvement over the default settings. Organelles were labelled based on subcellular  
753 locations assigned by UniProt<sup>62</sup> to these proteins, zoom regions were annotated manually  
754 based on available literature. Plot coordinates and annotations are available as  
755 Supplementary Table 4.

### 756 **Network visualizations**

757 In addition to t-SNE, the protein co-regulation matrix was also visualized as an undirected,  
758 weighted network using the igraph<sup>105</sup> and GGally<sup>106</sup> packages in R. The network contains the  
759 same 5,013 proteins as the co-regulation map, but only considers links above the arbitrary  
760 co-regulation threshold, i.e. between the top-scoring 0.5% of protein pairs. For these pairs,  
761 the network edges are weighted by the co-regulation score. A set of common network layout  
762 algorithms were deployed through the sna (social network analysis)<sup>107</sup> R package.

### 763 **Testing for co-functionality among of co-regulated proteins**

764 To test if protein co-regulation reflects co-function we defined three sets of “functionally  
765 related” protein pairs (subunits of the same protein complexes, enzymes catalyzing  
766 consecutive metabolic reactions and proteins with identical subcellular localization) as  
767 previously described<sup>25</sup>.

768 To test larger groups (not pairs) of co-regulated proteins for functional enrichment, we  
769 analyzed enrichment of Gene Ontology terms using the topGO<sup>108</sup> R package. For each  
770 protein we tested the group of its co-regulation partners for GO term enrichment. Because  
771 some proteins are co-regulated with no or very few other proteins, we restricted the analysis  
772 to proteins that are co-regulated with at least 10 proteins. The three aspects (Biological  
773 process, Molecular function, Cellular component) of GO were downloaded from QuickGO<sup>109</sup>  
774 with taxon set to human and qualifier to null. Rather than the whole proteome, only proteins  
775 that were included in the treeClust analysis and had GO annotations were used as the gene  
776 “universe” or background for the topGO analysis. Enrichment of GO terms among protein  
777 co-regulation groups was tested considering GO graph structure and using a Fisher’s exact  
778 test.

### 779 **Annotation of the co-regulation map**

780 Proteins localizing to specific subcellular compartments were downloaded from UniProt<sup>62</sup>  
781 using the following tags: Nucleus (SL-0191), Nucleolus (SL-0188), Endoplasmic reticulum  
782 (SL-0095), Mitochondrion (SL-0173), Cytoplasm (SL-0086), Secreted (SL-0243). Proteins  
783 and protein complexes in zoom regions (Fig. 1i) were annotated individually based on the  
784 available literature.

### 785 **Creating the www.proteomeHD.net framework**

786 The ProteomeHD online application was written in Python Flask web framework. The  
787 interactive plots are generated using Bokeh visualization library for Python

788 (<https://github.com/bokeh/bokeh>). The Gene Ontology and KEGG enrichment statistics are  
789 obtained from a STRING<sup>61</sup> server using an API call with maximally top 100 proteins  
790 co-regulated with the query. Only significantly enriched terms (hypergeometric test,  
791 Bonferroni adjusted  $P$  value  $< 0.1$ ) are displayed.

## 792 **Comparison to orthogonal methods**

793 Physical protein-protein-interactions (PPIs) detected by a comprehensive range of small-  
794 and large-scale methods were assessed using BioGRID<sup>60</sup>, version 3.4.152. Data from  
795 IntAct<sup>59</sup> were used as a smaller but curated resource of physical PPIs. Functional protein  
796 associations mapped by a large range of methods and publications were inferred from  
797 STRING<sup>61</sup>, version 10.5. Note that the protein co-regulation scores described here are only  
798 used by STRING starting with version 11<sup>76</sup>. BioPlex 2.0<sup>4</sup> served as an example for physical  
799 interactions mapped by a single project.

## 800 **Annotation of uncharacterized and disease genes**

801 Proteins were defined as “uncharacterized” on the basis of having an annotation score  $\leq 3$  in  
802 UniProt<sup>62</sup>. The UniProt annotation score is a heuristic measure of the annotation state of a  
803 protein, expressed as a 5-point system ([www.uniprot.org/help/annotation\\_score](http://www.uniprot.org/help/annotation_score)). The score  
804 combines various types and layers of UniProt annotation, and weights manually curated  
805 evidence higher than automated annotation. It may not always agree with the state of  
806 “characterization” that field experts would assign to the same protein. However, as an  
807 unbiased, data-driven approach we believe the UniProt annotation score is better suited to  
808 systematically identify uncharacterized proteins than manual annotation could be. Even with  
809 a systematic way of measuring the degree of annotation, the definition of what constitutes an  
810 “uncharacterised” protein is an arbitrary one. We chose “3 points or less” as the  
811 “uncharacterized” cut-off, because the available information for such proteins tends to be  
812 very vague, e.g. a sequence-based prediction as “multi-pass membrane protein”. In contrast,  
813 we found that the biological function of most 4-star proteins could be established reasonably  
814 well from the available literature.

815 The Cancer Gene Census, i.e. genes that can cause cancer when mutated, was  
816 curated by COSMIC (Catalogue Of Somatic Mutations In Cancer, version 81)<sup>63</sup>. DisGeNET  
817 was used as a comprehensive, curated list of human gene - disease associations<sup>64</sup>.

## 818 **Comparison of mRNA and protein expression profiling**

819 For the comparison of matched samples and proteins we considered mRNA and protein  
820 expression changes across 59 lymphoblastoid cell lines (Fig. 4a). The protein fold-changes  
821 are part of ProteomeHD and were originally published by Battle and colleagues<sup>30</sup>.  
822 RNA-sequencing data for the same cell lines and proteins were also previously reported<sup>86</sup>.  
823 We used the RNA-sequencing data to calculate mRNA fold-changes relative to a 60th cell  
824 line, which was the same cell line used as a SILAC reference for the protein expression data.  
825 The combined mRNA and protein dataset has been described in more detail elsewhere<sup>25</sup>.  
826 Fold-changes for genes covered by both the transcriptomics and proteomics analysis were  
827 subjected to treeClust learning (default parameters) and PR curves were obtained as  
828 described above.



829 For a more comprehensive comparison we considered protein associations predicted  
830 using treeClust learning or PCC on the basis of all 294 SILAC ratios in ProteomeHD (Fig.  
831 4b). This was compared to mRNA associations inferred by PCC on the basis of all human  
832 mRNA expression data processed by STRING. STRING's state-of-the-art mRNA  
833 coexpression analysis pipeline considers all microarray and RNA-sequencing data deposited  
834 in the GEO repository<sup>87</sup>, resulting in one of the largest mRNA coexpression analyses  
835 available to date<sup>61,88</sup>. Note that for this comparison we did not use the STRING coexpression  
836 score, which is calibrated against the KEGG database, but the original uncalibrated  
837 Pearson's correlations, which were kindly provided by Damian Szklarczyk. STRING PCCs  
838 are calculated separately for one- and two-channel microarrays and RNA-sequencing  
839 experiments. We used the average of these for the precision - recall analysis, which  
840 performed better than any individual experiment type.

#### 841 **Validation of treeClust and t-SNE on the cancer proteomics dataset**

842 Lapek *et al* measured the abundances for 6,911 proteins in 41 different breast cancer cell  
843 lines<sup>20</sup>. These data are available as Supplementary Table 2 (tab 3) of their report. As  
844 described by Lapek *et al*, we converted the protein intensities into log2 fold-changes over the  
845 median intensity measured for each protein across all cell lines. We then calculated  
846 Pearson's, Spearman's rank and bicor correlations for all possible protein pairs, as for  
847 ProteomeHD. The Spearman's correlation coefficients obtained in this way are identical to  
848 the ones obtained by Lapek *et al* using the cor.prob function (Supplementary Table 6 in their  
849 report<sup>20</sup>). We also determined treeClust co-regulation scores for all protein pairs. However,  
850 treeClust can only grow one decision tree per input variable, i.e. 41 in this dataset, which  
851 would be too few for it to perform properly. To circumvent this, we forced treeClust to  
852 generate 1,000 decision trees by applying it iteratively. We created 100 treeClust forests,  
853 each generated with a random subset of 10 of the 41 variables, and used the average  
854 co-regulation score for downstream analysis. Precision-recall analysis using a Reactome  
855 gold standard and t-SNE visualization were performed as described above. The CORUM  
856 protein complexes displayed in Lapek *et al*'s Figure 2, reported in their Supplementary Table  
857 7<sup>20</sup>, were color-coded in the co-regulation map.

#### 858 **Comparison of protein co-regulation and co-occurrence**

859 Two different approaches were used to measure protein co-occurrence in ProteomeHD.  
860 First, the Jaccard / Tanimoto similarity coefficient<sup>53</sup> was calculated using the Jaccard  
861 package for R. Second, a binary version of ProteomeHD was created, where all SILAC ratios  
862 were represented by 1s ("protein quantified") and all missing values were turned to 0s  
863 ("protein not quantified"). Subsequently, treeClust dissimilarities were re-calculated based on  
864 this binary version of ProteomeHD. The performance of these different metrics was  
865 assessed by a precision - recall analysis as described above.

#### 866 **Plasmids, siRNA, and antibodies**

867 For cloning of peroxisome-targeted Miro1, the C-terminal TMD and tail of Myc-Miro1 (kindly  
868 provided by P. Aspenström, Karolinska Institute, Sweden) was exchanged by a PEX26/ALDP  
869 fragment previously shown to target proteins to the peroxisome membrane<sup>82</sup>. PEX11β-EGFP

870 was kindly provided by G. Dodt (Univ. of Tuebingen, Germany). PEX11 $\beta$  siRNA (AUU AGG  
871 GUG AGA AUA GAC AGG AUGG) (Eurofins) was previously verified<sup>110</sup>. Control siRNA  
872 (si-GENOME nontargeting siRNA pool #2) was obtained from GE Healthcare  
873 (D-001206-14-05). Antibodies used were as follows: rabbit polyclonal antibody against  
874 PEX14 (1:1400, kindly provided by D. Crane, Griffith University, Australia); mouse  
875 monoclonal antibody 9E10 against the Myc epitope (1:200, Santa Cruz Biotechnology, Inc.,  
876 sc-40), rabbit monoclonal antibody against PEX11 $\beta$  (1:1000, Abcam, ab181066); rabbit  
877 polyclonal antibody against GAPDH (1:2000, ProSci3783). Secondary anti-IgG antibodies  
878 against rabbit (Alexa 594, 1:1000, Molec. Probes/Life Technol. A21207) and mouse (Alexa  
879 488, 1:400, Molec. Probes/Life Technol. A21202) were obtained from ThermoFisher  
880 Scientific. HRP-coupled donkey polyclonal antibody against rabbit IgG (1:5000) was  
881 obtained from Biorad (172-1013).

### 882 **Cell culture and transfection**

883 COS-7 cells (African green monkey kidney cells; ATCC CRL-1651), and PEX5 deficient  
884 fibroblasts (kindly provided by H. Waterham, AMC, University of Amsterdam, NL) were  
885 cultured in DMEM (high glucose, 4.5 g/L) supplemented with 10% FBS, 100 U/ml penicillin  
886 and 100  $\mu$ g/ml streptomycin at 37°C (5% CO<sub>2</sub>, 95% humidity) (HERACell 240i CO<sub>2</sub>  
887 incubator). COS-7 cells were transfected using diethylaminoethyl-dextran (Sigma-Aldrich).  
888 dPEX5 fibroblasts have enlarged peroxisomes, which facilitates the visualization of  
889 membrane extensions. For transfection of dPEX5 fibroblasts, the Neon® Transfection  
890 System (Thermo Fisher Scientific) was used following the manufacturer's protocol. Briefly,  
891 cells (seeded 24h before transfection) were washed once with PBS and trypsinized using  
892 TrypLE Express. Trypsinized cells were resuspended in complete medium, pelleted by  
893 centrifugation, and washed with PBS. The cells were once again centrifuged and carefully  
894 resuspended in 110  $\mu$ l buffer R. For each condition, 4  $\times$  10<sup>5</sup> cells were mixed with the DNA  
895 construct (5  $\mu$ g) or with 100 nM siRNA. Cells were microporated using a 100  $\mu$ l Neon tip with  
896 the following settings: 1400 V, 20 ms, one pulse. Microporated cells were immediately  
897 seeded into plates with prewarmed complete medium (without antibiotics) and incubated at  
898 37°C with 5% CO<sub>2</sub> and 95% humidity. The efficiency of silencing was monitored by  
899 immunoblotting of cell lysates and confirmed as previously reported<sup>110</sup>.

### 900 **Immunofluorescence and microscopy**

901 Cells grown on glass coverslips were processed for immunofluorescence 24h after  
902 transfection. Cells were fixed for 20 min with 4% paraformaldehyde in PBS (pH 7.4),  
903 permeabilized with 0.2% Triton X-100, and blocked with 1% BSA, each for 10 min.  
904 Incubation with primary and secondary antibodies took place for 1h each in a humid  
905 chamber. Coverslips were washed with ddH<sub>2</sub>O to remove PBS and mounted with Mowiol  
906 medium on glass slides. All immunofluorescence steps were performed at room temperature  
907 and cells were washed three times with PBS between each individual step. Cell imaging was  
908 performed using an IX81 microscope (Olympus) equipped with an UPlanSApo 100 $\times$ /1.40 oil  
909 objective (Olympus). Digital images were taken with a CoolSNAP HQ2 CCD camera and  
910 adjusted for contrast and brightness using the Olympus Soft Imaging Viewer software and  
911 MetaMorph 7 (Molecular Devices). For live-cell imaging, COS-7 cells were plated in 3.5 cm

912 diameter glass bottom dishes (Cellvis). MitoTracker Red CMXRos (Life Technologies) at 100  
913 nM was used for visualisation of mitochondria. Live-cell imaging data was collected using an  
914 Olympus IX81 microscope equipped with a Yokogawa CSUX1 spinning disk head,  
915 CoolSNAP HQ2 CCD camera, 60 x/1.35 oil objective. Digital images were taken and  
916 processed using VisiView software (Visitron Systems, Germany). Prior to image acquisition,  
917 a controlled temperature chamber was set-up on the microscope stage at 37°C, as well as  
918 an objective warmer. During image acquisition, cells were kept at 37°C and in  
919 CO<sub>2</sub>-independent medium (HEPES buffered). 200 stacks of 9 planes (0.5 µm thickness, 100  
920 ms exposure) were taken in a continuous stream. All conditions and laser intensities were  
921 kept between experiments.

#### 922 **Quantification and statistical analysis of peroxisome morphology and interaction**

923 Analysis of statistical significance was performed using GraphPad Prism 5 software. A  
924 two-tailed unpaired *t* test was used to determine statistical difference against the indicated  
925 group. \**P* < 0.05, \*\**P* < 0.01, \*\*\**P* < 0.001. For analysis of peroxisome morphology, a  
926 minimum of 150 cells were examined per condition, and organelle parameters (e.g.  
927 membrane protrusions) were microscopically assessed in at least three independent  
928 experiments. The analysis was made blind and in different areas of the coverslip. Organelle  
929 interaction and contact time were analysed manually from live-cell imaging data using  
930 MetaMorph 7 (Molecular Devices). A region of interest (ROI) was drawn in different areas of  
931 the cell. Spherical and elongated peroxisomes within the ROI were tracked over the whole  
932 time course, and the frequency and duration of contacts monitored. Multiple interactions of  
933 the same peroxisome with mitochondria were treated as separate events. Data are  
934 presented as mean ± SD.

935 **REFERENCES**

- 936 1. Gavin, A.-C. *et al.* Proteome survey reveals modularity of the yeast cell machinery.  
937 *Nature* **440**, 631–636 (2006).
- 938 2. Havugimana, P. C. *et al.* A census of human soluble protein complexes. *Cell* **150**,  
939 1068–1081 (2012).
- 940 3. Hein, M. Y. *et al.* A human interactome in three quantitative dimensions organized by  
941 stoichiometries and abundances. *Cell* **163**, 712–723 (2015).
- 942 4. Huttlin, E. L. *et al.* Architecture of the human interactome defines protein communities  
943 and disease networks. *Nature* **545**, 505–509 (2017).
- 944 5. Rolland, T. *et al.* A proteome-scale map of the human interactome network. *Cell* **159**,  
945 1212–1226 (2014).
- 946 6. Dunkley, T. P. J., Watson, R., Griffin, J. L., Dupree, P. & Lilley, K. S. Localization of  
947 organelle proteins by isotope tagging (LOPIT). *Mol. Cell. Proteomics* **3**, 1128–1134  
948 (2004).
- 949 7. Foster, L. J. *et al.* A mammalian organelle map by protein correlation profiling. *Cell* **125**,  
950 187–199 (2006).
- 951 8. Christoforou, A. *et al.* A draft map of the mouse pluripotent stem cell spatial proteome.  
952 *Nat. Commun.* **7**, 8992 (2016).
- 953 9. Thul, P. J. *et al.* A subcellular map of the human proteome. *Science* **356**, (2017).
- 954 10. Costanzo, M. *et al.* A global genetic interaction network maps a wiring diagram of  
955 cellular function. *Science* **353**, (2016).
- 956 11. Müllleder, M. *et al.* Functional Metabolomics Describes the Yeast Biosynthetic  
957 Regulome. *Cell* **167**, 553–565.e12 (2016).
- 958 12. Schena, M., Shalon, D., Davis, R. W. & Brown, P. O. Quantitative monitoring of gene  
959 expression patterns with a complementary DNA microarray. *Science* **270**, 467–470  
960 (1995).
- 961 13. DeRisi, J. L., Iyer, V. R. & Brown, P. O. Exploring the metabolic and genetic control of  
962 gene expression on a genomic scale. *Science* **278**, 680–686 (1997).
- 963 14. Eisen, M. B., Spellman, P. T., Brown, P. O. & Botstein, D. Cluster analysis and display of  
964 genome-wide expression patterns. *Proc. Natl. Acad. Sci. U. S. A.* **95**, 14863–14868  
965 (1998).
- 966 15. Kim, S. K. *et al.* A gene expression map for *Caenorhabditis elegans*. *Science* **293**,  
967 2087–2092 (2001).
- 968 16. Hughes, T. R. *et al.* Functional discovery via a compendium of expression profiles. *Cell*  
969 **102**, 109–126 (2000).
- 970 17. Stuart, J. M., Segal, E., Koller, D. & Kim, S. K. A gene-coexpression network for global  
971 discovery of conserved genetic modules. *Science* **302**, 249–255 (2003).
- 972 18. Singh, S. A. *et al.* Co-regulation proteomics reveals substrates and mechanisms of  
973 APC/C-dependent degradation. *EMBO J.* **33**, 385–399 (2014).
- 974 19. Wang, J. *et al.* Proteome Profiling Outperforms Transcriptome Profiling for Coexpression  
975 Based Gene Function Prediction. *Mol. Cell. Proteomics* **16**, 121–134 (2017).
- 976 20. Lapek, J. D., Jr *et al.* Detection of dysregulated protein-association networks by  
977 high-throughput proteomics predicts cancer vulnerabilities. *Nat. Biotechnol.* **35**, 983–989  
978 (2017).

- 979 21. Liu, Y., Beyer, A. & Aebersold, R. On the Dependency of Cellular Protein Levels on  
980 mRNA Abundance. *Cell* **165**, 535–550 (2016).
- 981 22. Wilhelm, M. *et al.* Mass-spectrometry-based draft of the human proteome. *Nature* **509**,  
982 582–587 (2014).
- 983 23. Fortelny, N., Overall, C. M., Pavlidis, P. & Freue, G. V. C. Can we predict protein from  
984 mRNA levels? *Nature* **547**, E19–E20 (2017).
- 985 24. Batada, N. N., Urrutia, A. O. & Hurst, L. D. Chromatin remodelling is a major source of  
986 coexpression of linked genes in yeast. *Trends Genet.* **23**, 480–484 (2007).
- 987 25. Kustatscher, G., Grabowski, P. & Rappsilber, J. Pervasive coexpression of spatially  
988 proximal genes is buffered at the protein level. *Mol. Syst. Biol.* **13**, 937 (2017).
- 989 26. Hurst, L. D. It's easier to get along with the quiet neighbours. *Mol. Syst. Biol.* **13**, 943  
990 (2017).
- 991 27. Raj, A., Peskin, C. S., Tranchina, D., Vargas, D. Y. & Tyagi, S. Stochastic mRNA  
992 synthesis in mammalian cells. *PLoS Biol.* **4**, e309 (2006).
- 993 28. Ebisuya, M., Yamamoto, T., Nakajima, M. & Nishida, E. Ripples from neighbouring  
994 transcription. *Nat. Cell Biol.* **10**, 1106–1113 (2008).
- 995 29. Khan, Z. *et al.* Primate transcript and protein expression levels evolve under  
996 compensatory selection pressures. *Science* **342**, 1100–1104 (2013).
- 997 30. Battle, A. *et al.* Genomic variation. Impact of regulatory variation from RNA to protein.  
998 *Science* **347**, 664–667 (2015).
- 999 31. Geiger, T., Cox, J. & Mann, M. Proteomic changes resulting from gene copy number  
1000 variations in cancer cells. *PLoS Genet.* **6**, e1001090 (2010).
- 1001 32. Stingele, S. *et al.* Global analysis of genome, transcriptome and proteome reveals the  
1002 response to aneuploidy in human cells. *Mol. Syst. Biol.* **8**, 608 (2012).
- 1003 33. Dephoure, N. *et al.* Quantitative proteomic analysis reveals posttranslational responses  
1004 to aneuploidy in yeast. *Elife* **3**, e03023 (2014).
- 1005 34. Ohta, S. *et al.* The protein composition of mitotic chromosomes determined using  
1006 multiclassifier combinatorial proteomics. *Cell* **142**, 810–821 (2010).
- 1007 35. Wu, L. *et al.* Variation and genetic control of protein abundance in humans. *Nature* **499**,  
1008 79–82 (2013).
- 1009 36. Kustatscher, G. *et al.* Proteomics of a fuzzy organelle: interphase chromatin. *EMBO J.*  
1010 **33**, 648–664 (2014).
- 1011 37. Wu, Y. *et al.* Multilayered genetic and omics dissection of mitochondrial activity in a  
1012 mouse reference population. *Cell* **158**, 1415–1430 (2014).
- 1013 38. Kustatscher, G., Grabowski, P. & Rappsilber, J. Multiclassifier combinatorial proteomics  
1014 of organelle shadows at the example of mitochondria in chromatin data. *Proteomics* **16**,  
1015 393–401 (2016).
- 1016 39. Okada, H., Ebhardt, H. A., Vonesch, S. C., Aebersold, R. & Hafen, E. Proteome-wide  
1017 association studies identify biochemical modules associated with a wing-size phenotype  
1018 in *Drosophila melanogaster*. *Nat. Commun.* **7**, 12649 (2016).
- 1019 40. Williams, E. G. *et al.* Systems proteomics of liver mitochondria function. *Science* **352**,  
1020 aad0189 (2016).
- 1021 41. Gupta, S., Turan, D., Tavernier, J. & Martens, L. The online Tabloid Proteome: an  
1022 annotated database of protein associations. *Nucleic Acids Res.* (2017).

- 1023 doi:10.1093/nar/gkx930
- 1024 42. Rieckmann, J. C. *et al.* Social network architecture of human immune cells unveiled by  
1025 quantitative proteomics. *Nat. Immunol.* **18**, 583–593 (2017).
- 1026 43. Kim, M.-S. *et al.* A draft map of the human proteome. *Nature* **509**, 575–581 (2014).
- 1027 44. Uhlén, M. *et al.* Proteomics. Tissue-based map of the human proteome. *Science* **347**,  
1028 1260419 (2015).
- 1029 45. Ong, S.-E. *et al.* Stable isotope labeling by amino acids in cell culture, SILAC, as a  
1030 simple and accurate approach to expression proteomics. *Mol. Cell. Proteomics* **1**,  
1031 376–386 (2002).
- 1032 46. Vizcaíno, J. A. *et al.* 2016 update of the PRIDE database and its related tools. *Nucleic*  
1033 *Acids Res.* **44**, D447–56 (2016).
- 1034 47. Fabregat, A. *et al.* The Reactome pathway Knowledgebase. *Nucleic Acids Res.* **44**,  
1035 D481–7 (2016).
- 1036 48. Buttrely, S. E. & Whitaker, L. R. treeClust: an R package for tree-based clustering  
1037 dissimilarities. *The R Journal* **7**, 227–236 (2015).
- 1038 49. Buttrely, S. E. & Whitaker, L. R. A scale-independent, noise-resistant dissimilarity for  
1039 tree-based clustering of mixed data. *NPS Technical Report Archive* (2016). Available at:  
1040 <https://calhoun.nps.edu/handle/10945/48615>.
- 1041 50. Ravasz, E., Somera, A. L., Mongru, D. A., Oltvai, Z. N. & Barabási, A. L. Hierarchical  
1042 organization of modularity in metabolic networks. *Science* **297**, 1551–1555 (2002).
- 1043 51. Yip, A. M. & Horvath, S. Gene network interconnectedness and the generalized  
1044 topological overlap measure. *BMC Bioinformatics* **8**, 22 (2007).
- 1045 52. The Gene Ontology Consortium. Expansion of the Gene Ontology knowledgebase and  
1046 resources. *Nucleic Acids Res.* **45**, D331–D338 (2017).
- 1047 53. Jaccard, P. Distribution de la flore alpine dans le bassin des Dranses et dans quelques  
1048 régions voisines. *Bull. Soc. Vaud. sci. nat.* **37**, 241–272 (1901).
- 1049 54. Kustatscher, G., Grabowski, P. & Rappsilber, J. treeClust improves protein co-regulation  
1050 analysis due to robust selectivity for close linear relationships. *bioRxiv* (2019).  
1051 doi:10.1101/578971
- 1052 55. Krzywinski, M., Birol, I., Jones, S. J. M. & Marra, M. A. Hive plots—rational approach to  
1053 visualizing networks. *Brief. Bioinform.* **13**, 627–644 (2012).
- 1054 56. Van Der Maaten, L. & Hinton, G. Visualizing High-Dimensional Data Using t-SNE. *J.*  
1055 *Mach. Learn. Res.* **9**, 26 (2008).
- 1056 57. Van Der Maaten, L. Accelerating t-SNE using tree-based algorithms. *J. Mach. Learn.*  
1057 *Res.* **15**, 3221–3245 (2014).
- 1058 58. García-Aguilar, A. & Cuezva, J. M. A Review of the Inhibition of the Mitochondrial ATP  
1059 Synthase by IF1 in vivo: Reprogramming Energy Metabolism and Inducing  
1060 Mitohormesis. *Front. Physiol.* **9**, 1322 (2018).
- 1061 59. Orchard, S. *et al.* The MIntAct project—IntAct as a common curation platform for 11  
1062 molecular interaction databases. *Nucleic Acids Res.* **42**, D358–63 (2014).
- 1063 60. Stark, C. *et al.* BioGRID: a general repository for interaction datasets. *Nucleic Acids*  
1064 *Res.* **34**, D535–9 (2006).
- 1065 61. Szklarczyk, D. *et al.* The STRING database in 2017: quality-controlled protein-protein  
1066 association networks, made broadly accessible. *Nucleic Acids Res.* **45**, D362–D368

- 1067 (2017).
- 1068 62. The UniProt Consortium. UniProt: the universal protein knowledgebase. *Nucleic Acids*  
1069 *Res.* **45**, D158–D169 (2017).
- 1070 63. Forbes, S. A. *et al.* COSMIC: somatic cancer genetics at high-resolution. *Nucleic Acids*  
1071 *Res.* **45**, D777–D783 (2017).
- 1072 64. Piñero, J. *et al.* DisGeNET: a comprehensive platform integrating information on human  
1073 disease-associated genes and variants. *Nucleic Acids Res.* **45**, D833–D839 (2017).
- 1074 65. Andrews, S. J. & Rothnagel, J. A. Emerging evidence for functional peptides encoded  
1075 by short open reading frames. *Nat. Rev. Genet.* **15**, 193–204 (2014).
- 1076 66. Aspden, J. L. *et al.* Extensive translation of small Open Reading Frames revealed by  
1077 Poly-Ribo-Seq. *Elife* **3**, e03528 (2014).
- 1078 67. D’Lima, N. G. *et al.* A human microprotein that interacts with the mRNA decapping  
1079 complex. *Nat. Chem. Biol.* **13**, 174–180 (2017).
- 1080 68. Chu, Q. *et al.* Identification of Microprotein-Protein Interactions via APEX Tagging.  
1081 *Biochemistry* (2017). doi:10.1021/acs.biochem.7b00265
- 1082 69. Slavoff, S. A. *et al.* Peptidomic discovery of short open reading frame-encoded peptides  
1083 in human cells. *Nat. Chem. Biol.* **9**, 59–64 (2013).
- 1084 70. Meyer, B., Wittig, I., Trifilieff, E., Karas, M. & Schägger, H. Identification of two proteins  
1085 associated with mammalian ATP synthase. *Mol. Cell. Proteomics* **6**, 1690–1699 (2007).
- 1086 71. Chen, R., Runswick, M. J., Carroll, J., Fearnley, I. M. & Walker, J. E. Association of two  
1087 proteolipids of unknown function with ATP synthase from bovine heart mitochondria.  
1088 *FEBS Lett.* **581**, 3145–3148 (2007).
- 1089 72. Fujikawa, M., Ohsakaya, S., Sugawara, K. & Yoshida, M. Population of ATP synthase  
1090 molecules in mitochondria is limited by available 6.8-kDa proteolipid protein (MLQ).  
1091 *Genes Cells* **19**, 153–160 (2014).
- 1092 73. Borner, G. H. H. *et al.* Multivariate proteomic profiling identifies novel accessory proteins  
1093 of coated vesicles. *J. Cell Biol.* **197**, 141–160 (2012).
- 1094 74. Signorile, A., Sgaramella, G., Bellomo, F. & De Rasmio, D. Prohibitins: A Critical Role in  
1095 Mitochondrial Functions and Implication in Diseases. *Cells* **8**, (2019).
- 1096 75. Brennan, R. *et al.* Investigating nucleo-cytoplasmic shuttling of the human DEAD-box  
1097 helicase DDX3. *Eur. J. Cell Biol.* **97**, 501–511 (2018).
- 1098 76. Szklarczyk, D. *et al.* STRING v11: protein-protein association networks with increased  
1099 coverage, supporting functional discovery in genome-wide experimental datasets.  
1100 *Nucleic Acids Res.* **47**, D607–D613 (2019).
- 1101 77. Schrader, M., Costello, J. L., Godinho, L. F., Azadi, A. S. & Islinger, M. Proliferation and  
1102 fission of peroxisomes - An update. *Biochim. Biophys. Acta* **1863**, 971–983 (2016).
- 1103 78. Schrader, M., Costello, J., Godinho, L. F. & Islinger, M. Peroxisome-mitochondria  
1104 interplay and disease. *J. Inherit. Metab. Dis.* **38**, 681–702 (2015).
- 1105 79. Devine, M. J., Birsa, N. & Kittler, J. T. Miro sculpts mitochondrial dynamics in neuronal  
1106 health and disease. *Neurobiol. Dis.* **90**, 27–34 (2016).
- 1107 80. Costello, J. L. *et al.* Predicting the targeting of tail-anchored proteins to subcellular  
1108 compartments in mammalian cells. *J. Cell Sci.* **130**, 1675–1687 (2017).
- 1109 81. Okumoto, K. *et al.* New splicing variants of mitochondrial Rho GTPase-1 (Miro1)  
1110 transport peroxisomes. *J. Cell Biol.* **217**, 619–633 (2018).

- 1111 82. Castro, I. G. *et al.* A role for Mitochondrial Rho GTPase 1 (MIRO1) in motility and  
1112 membrane dynamics of peroxisomes. *Traffic* **19**, 229–242 (2018).
- 1113 83. Rodríguez-Serrano, M., Romero-Puertas, M. C., Sanz-Fernández, M., Hu, J. &  
1114 Sandalio, L. M. Peroxisomes Extend Peroxules in a Fast Response to Stress via a  
1115 Reactive Oxygen Species-Mediated Induction of the Peroxin PEX11a. *Plant Physiol.*  
1116 **171**, 1665–1674 (2016).
- 1117 84. Mattiazzi Ušaj, M. *et al.* Genome-Wide Localization Study of Yeast Pex11 Identifies  
1118 Peroxisome-Mitochondria Interactions through the ERMES Complex. *J. Mol. Biol.* **427**,  
1119 2072–2087 (2015).
- 1120 85. Shai, N. *et al.* Systematic mapping of contact sites reveals tethers and a function for the  
1121 peroxisome-mitochondria contact. *Nat. Commun.* **9**, 1761 (2018).
- 1122 86. Pickrell, J. K. *et al.* Understanding mechanisms underlying human gene expression  
1123 variation with RNA sequencing. *Nature* **464**, 768–772 (2010).
- 1124 87. Barrett, T. *et al.* NCBI GEO: archive for functional genomics data sets--update. *Nucleic*  
1125 *Acids Res.* **41**, D991–5 (2013).
- 1126 88. Szklarczyk, D. *et al.* STRING v10: protein-protein interaction networks, integrated over  
1127 the tree of life. *Nucleic Acids Res.* **43**, D447–52 (2015).
- 1128 89. Gandhi, S. J., Zenklusen, D., Lionnet, T. & Singer, R. H. Transcription of functionally  
1129 related constitutive genes is not coordinated. *Nat. Struct. Mol. Biol.* **18**, 27–34 (2011).
- 1130 90. Jovanovic, M. *et al.* Immunogenetics. Dynamic profiling of the protein life cycle in  
1131 response to pathogens. *Science* **347**, 1259038 (2015).
- 1132 91. R Core Team. R: A Language and Environment for Statistical Computing. (2018).
- 1133 92. Dowle, M. & Srinivasan, A. data.table: Extension of `data.frame`. (2018).
- 1134 93. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer, 2016).
- 1135 94. Auguie, B. gridExtra: Miscellaneous Functions for 'Grid' Graphics. (2017).
- 1136 95. Wilke, C. O. cowplot: Streamlined Plot Theme and Plot Annotations for 'ggplot2'. (2018).
- 1137 96. Garnier, S. viridis: Default Color Maps from 'matplotlib'. (2018).
- 1138 97. Kustatscher, G., Wills, K. L. H., Furlan, C. & Rappsilber, J. Chromatin enrichment for  
1139 proteomics. *Nat. Protoc.* **9**, 2090–2099 (2014).
- 1140 98. Alabert, C. *et al.* Nascent chromatin capture proteomics determines chromatin dynamics  
1141 during DNA replication and identifies unknown fork components. *Nat. Cell Biol.* **16**,  
1142 281–293 (2014).
- 1143 99. Cox, J. & Mann, M. MaxQuant enables high peptide identification rates, individualized  
1144 p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nat. Biotechnol.*  
1145 **26**, 1367–1372 (2008).
- 1146 100. Zhang, B. & Horvath, S. A general framework for weighted gene co-expression network  
1147 analysis. *Stat. Appl. Genet. Mol. Biol.* **4**, Article17 (2005).
- 1148 101. Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network  
1149 analysis. *BMC Bioinformatics* **9**, 559 (2008).
- 1150 102. Langfelder, P. & Horvath, S. Fast R Functions for Robust Correlations and Hierarchical  
1151 Clustering. *J. Stat. Softw.* **46**, (2012).
- 1152 103. Sing, T., Sander, O., Beerenwinkel, N. & Lengauer, T. ROCr: visualizing classifier  
1153 performance in R. *Bioinformatics* **21**, 3940–3941 (2005).
- 1154 104. Krijthe, J. H. Rtsne: T-Distributed Stochastic Neighbor Embedding using Barnes-Hut



- 1155 Implementation. URL: <https://github.com/jkrijthe/Rtsne> (2015).
- 1156 105.Csardi, G. & Nepusz, T. The igraph software package for complex network research.  
1157 *InterJournal*, 1695 (2006).
- 1158 106.Schloerke, B. *et al.* GGally: Extension to 'ggplot2'. (2018).
- 1159 107.Butts, C. T. sna: Tools for Social Network Analysis. (2016).
- 1160 108.Alexa, A. & Rahnenfuhrer, J. topGO: enrichment analysis for gene ontology. *R package*  
1161 *version 2.30.0* (2016).
- 1162 109.Binns, D. *et al.* QuickGO: a web-based tool for Gene Ontology searching. *Bioinformatics*  
1163 **25**, 3045–3046 (2009).
- 1164 110.Costello, J. L. *et al.* ACBD5 and VAPB mediate membrane associations between  
1165 peroxisomes and the ER. *J. Cell Biol.* **216**, 331–342 (2017).