

1 **Methodological insights on proteogenomic approaches** 2 **to enhance proteomics**

3 Laura Fancello¹ and Thomas Burger¹

4 ¹*Univ. Grenoble Alpes, CNRS, CEA, Inserm, Profi FR2048, Grenoble, France.*

5

6 **Abstract**

7 Proteogenomics aims at identifying variant or unknown proteins in bottom-up proteomics,
8 searching transcriptome- or genome-derived custom protein databases. However,
9 empirical observations reported that the large size of these proteogenomic databases is
10 associated to lower sensitivity of peptide identifications. Various strategies were proposed
11 to avoid this, including the generation of reduced transcriptome-informed protein
12 databases (*i.e.*, built from reference protein databases only retaining proteins with
13 expressed transcript in the sample-matched transcriptome), which were found to increase
14 peptide identification sensitivity. In this work, we propose a detailed evaluation of this
15 approach. First, we establish that the increased sensitivity in peptide identification is in fact
16 a statistical artefact, which directly results from the limited capability of TDC to accurately
17 model incorrect target matches with excessively small databases. As anti-conservative
18 FDRs likely hamper the robustness of the resulting biological conclusions, we advocate for
19 alternative FDR control methods that are less sensitive to database size. Second, we
20 show that despite not increasing sensitivity, reduced transcriptome-informed databases
21 are useful, as they allow reducing ambiguity of protein identifications, yielding fewer
22 shared peptides. Furthermore, we illustrate that searching the reference database and
23 subsequently filtering proteins with unexpressed transcript similarly reduces protein

24 identification ambiguity, while representing a more transparent and reproducible strategy.
25 To summarize, using reduced transcriptome-informed databases is an interesting strategy
26 that has not been promoted for the good reason (an artifactual peptide identification
27 sensitivity increment instead of a protein identification ambiguity decrement).

28

29 **Keywords:** proteogenomics, proteomics, transcriptome-informed protein databases,
30 peptide identification sensitivity, protein identification ambiguity, FDR control, target-decoy
31 competition

32

33 **BACKGROUND**

34 The term “proteogenomics” originally referred to the use of proteomics to enhance
35 genome annotation, by inferring coding genomic regions based on evidence from mass
36 spectrometry-based proteomics¹. Nowadays, it more broadly indicates the combined
37 analysis of genomics and/or transcriptomics together with proteomics in a large spectrum
38 of applications: from the study of gene expression regulation at transcript and protein level,
39 to the identification of specific protein variants expressed in cancer^{2,3}. Most importantly,
40 proteogenomics represents an attractive strategy to enhance proteomics in two main
41 ways: *i.* improving protein inference; *ii.* improving database searches for peptide
42 identification.

43 Protein inference is a central issue in proteomics, given the presence of shared
44 peptides. This nominates peptides that might originate from different proteins sharing
45 homology; or from different proteoforms due to alternative mRNA splicing, post-
46 translational modifications, proteolytic cleavages, and/or allelic variants. Indeed, in bottom-
47 up mass spectrometry-based proteomics, the most widely used proteomic approach,

48 peptide-protein connectivity is lost for experimental reasons and protein identifications are
49 to be inferred from peptide identifications. Traditionally, the issue of protein inference was
50 addressed using simple heuristics, such as the two-peptide rule (only proteins identified by
51 at least two peptides are retained) or the parsimonious principle (the smallest subset of
52 proteins which can explain most or all peptides is retained)^{4,5}. Later on, more refined
53 probability-based approaches were developed, which model shared peptide mappings to
54 their parent proteins⁶⁻⁹. Most commonly, when proteins cannot be discriminated based on
55 peptide identifications (*i.e.*, they are identified by the same set of peptides) they are
56 reported as a protein group, which complicates comparisons between different
57 experiments and protein quantification. In this context, proteogenomics can aid protein
58 inference using evidence from transcript expression: in particular, some Bayesian
59 approaches were developed based on this strategy¹⁰⁻¹². The other main contribution of
60 proteogenomics to proteomics relates to the refinement of reference protein databases
61 used for peptide identification. Classically, peptides resulting from bottom-up analyses are
62 identified by matching their experimentally measured mass spectra against theoretical
63 spectra of all candidate peptides from a user-selected reference protein database. The
64 underlying assumption is that such a database exhaustively and accurately describes all
65 protein sequences present in the sample. However, this may be unrealistic for two
66 reasons. First, reference databases only contain canonical -- experimentally validated or
67 predicted -- protein sequences, while other variants or isoforms may be present, especially
68 in tumour samples. Second, a reference protein database may simply be lacking for less
69 studied organisms with scarce or no genomic annotations. In the first case, more
70 exhaustive protein databases including undocumented or variant peptides can be
71 generated by appending to the reference database variant sequences from public genomic
72 repositories (*i.e.*, COSMIC or dbSNP)¹³⁻¹⁵ or sample-specific variants identified from

73 matched transcriptomes or genomes. In the second case, protein databases can be
74 generated by 6-frame translation of the genome or of the sample-matched
75 transcriptome^{16,17}. A major downfall of these proteogenomic databases is represented by
76 their typically large size. Searching very large databases represents a considerable
77 computational load and complicates the task of discriminating between correct and
78 incorrect matches. In particular, various works showed that when using target-decoy
79 competition (TDC) for FDR control on large database searches, fewer peptides are
80 identified at the same FDR level, which stands in stark contrast with the initial motivation of
81 proteogenomics¹⁸⁻²⁰. To avoid this, it was proposed to perform separate FDR validation of
82 canonical and novel peptides and to apply post-search filters or machine learning methods
83 to increase confidence in the newly identified peptides^{13-15,21,22}. Additionally, various
84 strategies were adopted to limit the size of databases generated by proteogenomics.
85 When possible, genome 6-frame translation was replaced by translating candidate ORFs
86 identified by gene prediction algorithms; sample-specific variants from matched
87 sequencing were preferentially added to the reference database rather than variant
88 sequences from COSMIC or dbSNP²³⁻²⁵; in some studies, after appending variants
89 identified from sequencing data, the reference database was reduced to only proteins with
90 transcript expressed according to transcriptomics, since according to the “central dogma of
91 biology”, there can be no protein without corresponding transcript²⁶⁻²⁸. It was also proposed
92 to generate reduced transcriptome-informed protein databases by barely reducing the
93 reference database to proteins with expressed transcript, without including any novel
94 sequence, with the only declared objective of increasing sensitivity in the identification of
95 known sequences^{26,28,29}. These works claim that searching such reduced transcriptome-
96 informed databases allows increasing the number of valid identifications. A strategy was
97 also proposed to optimize the balance between lost identifications due to the

98 incompleteness of an excessively reduced database and additional identifications
99 observed from searching reduced transcriptome-informed databases, so as to maximize
100 the number of valid identifications³⁰. However: *i.* Only limited attention was given to the
101 mechanistic explanation of the increased number of identified peptides with respect to
102 database size; *ii.* Little is known about the impact of reduced transcriptome-informed
103 database searches on protein inference, in terms of ambiguity of protein identification and
104 shared peptide assignments. Therefore, in this work, we investigated the use of reduced
105 transcriptome-informed sample-specific protein databases, focusing on these two
106 methodological aspects. Our investigations result into three conclusions. First, the reported
107 increment of the number of identifications obtained searching reduced transcriptome-
108 informed databases is a statistical artifact: it is only the spurious consequence of an
109 underestimated FDR, which results from the TDC sensitivity to the database size (also
110 reported in ³¹). In other words, reducing the search database to increase sensitivity broadly
111 amounts to validate peptide identifications at an FDR that is larger than reported, hereby
112 questioning the validity of peptides only identified thanks to a reduced database search
113 and comparability between studies. Second, searching reduced transcriptome-informed
114 protein databases followed by accurate FDR control remains nonetheless of interest, for it
115 decreases ambiguity of protein identifications by reducing the proportion of shared
116 peptides and the size of protein groups. Finally, searches against the reference database
117 followed by post-hoc filtering of proteins with no evidence of transcript expression provides
118 comparable proteomic identifications to searches against the reduced transcriptome-based
119 database, while guaranteeing more transparency and comparability between studies. We
120 therefore provided an R code to perform post-hoc filtering on proteomic identifications
121 based on transcript evidence and to manually inspect protein identifications and shared

122 peptide assignments within protein groups of interest together with transcript expression
123 information.

124

125

RESULTS

126 **Reduced transcriptome-informed database search does not increase**
127 **sensitivity if FDR is accurately controlled**

128 To investigate the impact of reduced transcriptome-informed protein databases on
129 proteomic identifications, we used two human samples (hereafter referred to as Jurkat and
130 Lung) for which matched transcriptome and proteome were publicly available (**Supp.**
131 **Table 1**). For each of them, we built a sample-specific reduced transcriptome-informed
132 protein database, in the following way. First, we processed transcriptome datasets and
133 identified the set of transcripts expressed in each sample using StringTie, a common
134 transcriptome assembly method (see “Transcriptome analysis” in Methods). Then, we
135 generated reduced databases for MS/MS search by retaining from the Ensembl human
136 protein database only those proteins whose transcript was expressed in the sample-
137 matched transcriptome (see “Construction of reduced transcriptome-informed protein
138 databases for MS/MS search” in Methods) (**Figure 1A-B**). We compared valid peptide-
139 spectrum matches (PSMs) obtained from the MS/MS search against the Ensembl human
140 database (referred to as the “full database”) or against the sample-specific reduced
141 database (referred to as the “reduced database”) at 1% FDR, as estimated by TDC. In
142 agreement with previous studies^{26,28,30}, we found that a few spectra and peptides identified
143 in the full database were lost in the reduced database search (“lost in reduced DB”), while
144 others were only identified in the reduced database search (“additional in reduced DB”)

145 **(Figure 1C)**. Lost identifications originate from reduced database incompleteness. Indeed,
146 even more identifications were lost when using a further reduced protein database, as the
147 one generated on the basis of the smaller set of expressed transcripts identified by
148 Cufflinks, an alternative method of transcriptome assembly (**Supp. Fig. 1**). Additional
149 identifications, instead, are commonly attributed to an increased sensitivity of MS/MS
150 searches against smaller databases, such as reduced transcriptome-informed protein
151 databases^{26,28,30}. In this work, we investigated more thoroughly the origin of additional
152 identifications obtained from searching these reduced transcriptome-informed protein
153 databases. To this end, we performed a detailed comparison of all (target or decoy) PSMs
154 retained from the full and reduced database searches, after validation prefilters (*i.e.*, single
155 best scoring PSM per spectrum, minimum peptide length of 7 amino acids), but prior to
156 filtering for 1% FDR control (**Figure 1A**). Since we built the reduced database as a simple
157 subset of the Ensembl human database and considered only a single best scoring peptide
158 per spectrum (see “Proteome analysis” in Methods), we could easily map each spectrum
159 match between the two searches. Two interesting observations emerged from this
160 comparison. First, several spectra are reallocated in the reduced database (*i.e.*, assigned
161 in the reduced database search to a different match from that of the full database), which
162 occurs when the peptide match from the full database is not included in the reduced
163 database. However, the PSM score is never higher in the reduced database: at best, it is
164 equal (data points on the diagonal, **Figure 2A, Supp. Fig. 2A**) to the score in the full
165 database, but for the most part it is smaller (**Figure 2B, Supp. Fig. 2B**). Therefore,
166 additional identifications in the reduced database do not come from an improved search
167 score. For sake of clarity, those few spectra with a match only in the reduced database
168 (indicated as “no match, target” or “no match, decoy” in **Figure 2A, Supp. Fig. 2A, 3A, 4**
169 **and 5**) do not contradict this observation; they are all explained by reallocation and

170 prefilters used for validation (**Supp. Fig. 3B**, see **Supplementary Note 1** for a detailed
171 explanation). Similarly, we also verified all cases of pretty rank (PSMs with score difference
172 > 0.1 , which are considered of equal score, see “Proteome analysis” section in Methods)
173 and again confirmed that the maximum PSM score in the reduced database in no case is
174 higher than in the full database search (**Supp. Fig. 3C-D**, see also **Supplementary Note**
175 **1**). The second main observation was that the score cutoff estimated by TDC at 1% FDR
176 (*i.e.*, the score defining the set of accepted PSMs, while respecting the constraint that less
177 than 1% of them are expected to be a false discoveries) is lower for the reduced database
178 than for the full database search (**Figure 2A**, **Supp. Fig. 2A**). Consequently, for a few
179 spectra their match is not validated after FDR control in the full database while, at lower or
180 equal score at best, it is validated in the reduced database search (pointed out by the
181 arrow in **Figure 2A** and **Supp. Fig. 2A**). This is clearly the reason why these PSMs are
182 accounted for as additional identifications in the reduced database search. We also
183 observed a few reallocations, which can likewise yield additional spectra and/or peptide
184 identifications in the reduced database search. They are, in particular, reallocations from
185 non-target matches in the full database to target matches in the reduced database search
186 (2.9% and 1.9% of all spectra in Jurkat and lung respectively) and reallocations between
187 different target matches (3.2% and 2.7% of all spectra in Jurkat and lung respectively)
188 (**Figure 2A,C** and **Supp. Fig. 2A,C**). However, only a minority of them are valid
189 identifications at 1% FDR control (*i.e.*, pass the score cutoff for FDR control from the
190 reduced database search) (**Figure 2C**, **Supp. Fig. 2C**, **Supp. Table 4**). Further, even
191 fewer of them would pass the cutoff obtained from the full database search and they are
192 hereafter named as “pure reallocations” to indicate that additional identifications from
193 these PSMs uniquely originate from reallocation and not from additionally validated PSMs
194 due to the the lower cutoff at 1% FDR validation in the reduced database (**Figure 2C**,

195 **Supp. Fig. 2C and 6A, Supp. Table 4**). Additional peptide identifications originating from
196 either lower cutoff for FDR control or from pure reallocations present inferior PSM scores
197 compared to peptide identifications obtained from both database searches (**Supp. Fig.**
198 **6B**). For pure reallocations, the difference in score between the full and reduced database
199 match can be quite important, especially for target PSMs in the full database, reallocated
200 in the reduced database search (**Supp. Fig. 6C**). Furthermore, additional peptide
201 identifications only allow obtaining 6 and 8 additional protein identifications (*i.e.*, protein
202 groups whose protein members are not identified in the full database search) in the Jurkat
203 and lung sample, respectively (**Supp. Table 5**). Thus, these additional identifications are of
204 lower quality and provide little benefit to protein identification. Overall, only few additional
205 identifications come from pure reallocations, while the main origin is the lower cutoff for
206 FDR control in the reduced database, explaining 98.6% (n=3,147) and 95.2% (n=1,875)
207 additional spectral identifications and 96.5% (n=5,560) and 77.5% (n=1,524) additional
208 peptide identifications for the Jurkat and lung samples, respectively (**Figure 2D-E, Supp.**
209 **Fig. 2D-E**). Therefore, we investigated the reasons why lower cutoffs were observed at a
210 same FDR threshold in the reduced databases. To do so, we first simulated what would
211 occur if it were instead equal to that of the full database (**Figure 3A, Supp. Fig. 7A**). We
212 observed that the proportion of valid decoys in the reduced database search would
213 considerably decrease compared to the full database, with a net loss of 38.4% and 27.1%
214 of valid decoys in the Jurkat (**Figure 3B**) and lung sample (**Supp. Fig. 7B**), respectively.
215 Indeed, an important fraction of spectra matching valid decoys in the full database are
216 assigned to invalid or non-decoy matches in the reduced database and not
217 counterbalanced by reallocations in the other direction (*i.e.*, from invalid/non-decoy
218 matches to valid decoys) (**Figure 3C, Supp. Fig. 7C, 4 and 5**). On the contrary, the
219 majority of spectra matching valid targets in the full database match the same valid target

220 in the reduced database, so that their loss is quite limited (**Figure 3C, Supp. Fig. 7C**).

221 While spectra matching valid decoys in the full database are reallocated much more

222 frequently in the reduced database than spectra matching valid targets, upon reallocation

223 they behave similarly: only few reallocations result in a valid match of the same type

224 (**Figure 3C, Supp. Fig. 7C**) and the score difference between full and reduced database

225 matches is comparable (**Supp. Fig. 8A**). Hence, the proportion of valid decoys lost in the

226 reduced database is higher than that of targets, simply because a higher proportion of

227 them is reallocated. This is easily explained by how the reduced database is generated:

228 only proteins whose transcript is expressed, thus those more likely to be present, are

229 retained from the canonical full protein database. Therefore, all valid targets from the full

230 database are in theory still present in the reduced database, while this is not the case for

231 decoys that represent by definition random matches (**Supp. Fig. 8B**). The lower cutoff

232 obtained by TDC for the reduced database allows to validate a few more decoys and thus

233 recover the proportion of valid decoys required to declare a nominal FDR level of 1%

234 (**Figure 3B, Supp. Fig. 7B**). We claim that additional identifications validated using a

235 lower cutoff in the reduced database represent a byproduct of the known influence of the

236 database size on TDC³¹, rather than an effect of increased sensitivity in reduced database

237 searches. Naturally, in absence of a benchmark, it is impossible to determine whether they

238 represent correct matches, missed in the full database due to FDR overestimation, or

239 incorrect matches, accepted in the reduced database due to FDR underestimation.

240 However, three main observations indicate that they should be at least considered with

241 caution. First, they are accepted in the reduced database at quite low scores, meaning

242 that, in any case, they represent low quality spectra and cannot be identified with very high

243 confidence. Second, it could be assumed that additional identifications stem from the

244 removal, in the reduced database, of high-scoring decoys out-competing correct target

245 matches, thus lowering sensitivity. However, in our study, most additional identifications do
246 not represent reallocations from decoys to targets; they consist, instead, in the same
247 PSMs, accepted at 1% FDR only in the reduced database because of a lower score cutoff
248 **(Figure 2D-E; Supp. Fig. 2D-E)**. Third and most importantly, the artifactual origin of
249 additional identifications, given TDC sensitivity to database size, is supported by the
250 comparison between the behavior of TDC and of the Benjamini-Hochberg (BH) procedure
251 for FDR control³². BH is known to be a conservative and stable FDR control procedure,
252 and it has recently been successfully applied to peptide identification³¹. In particular, TDC
253 was found to be less conservative and less stable than BH with respect to preliminary
254 filters on precursor mass accuracy: at narrower mass tolerance, fewer decoys were fair
255 competitors for incorrect random matches, artificially lowering cutoffs. Therefore, reducing
256 database size can similarly result into an insufficient number of decoys to accurately
257 simulate incorrect target matches and lead to the observed lower cutoffs. To confirm this,
258 we applied the BH procedure on target-only database searches (see “Proteome analysis”
259 section in Methods) and obtained more conservative score cutoffs and, most importantly,
260 more stable with respect to database size, compared to TDC **(Figure 3D, Supp. Fig. 7D,**
261 **Supp. Fig. 9A-B)**. Consistently, a much more limited number of additional identifications
262 was validated in reduced database searches using BH-based FDR control **(Supp. Fig.**
263 **9C)**. We also employed the BH procedure for FDR control on concatenated target-decoy
264 database searches; while doing so is a nonsense from a practical data processing
265 viewpoint, yet from a statistical methodology viewpoint, it simplifies comparative
266 evaluations of BH and TDC stabilities. As expected, we obtained more conservative and
267 stable score cutoffs with BH **(Supp. Fig. 9A-B)**.

268

269 **Transcriptome information aids to reduce ambiguity of protein identifications**

270 While not enhancing sensitivity of peptide identifications, reduced transcriptome-
271 informed databases can still benefit proteomics at the protein inference step, by lowering
272 ambiguity in protein identifications. These databases include fewer proteins – only those
273 proteins which are most likely present given their transcript expression – and it is
274 reasonable to assume that with fewer possible protein matches, we may obtain fewer
275 shared peptides and smaller protein groups. This decrement in protein group size has
276 already been observed, but either not discussed³³ or attributed to an additional number of
277 identifiable peptides available for parsimony-based protein inference²⁶. We already
278 illustrated how additional identifications from searching reduced databases actually
279 represent a flaw of TDC with respect to reduced database size, and how they can be
280 largely avoided using alternative procedures for FDR control, such as, for instance, BH.
281 We will now show that plain searches against reduced databases followed by BH-based
282 FDR control nonetheless yield smaller protein groups and less ambiguous protein
283 identifications, thus regardless of additional identifications or protein inference methods.
284 Concretely, we compared identifications obtained from the full or reduced database
285 searches followed by BH-based FDR control. The total number of identifications, at the
286 spectrum, peptide and protein level, is comparable (**Figure 4A**). As number of protein-level
287 identifications, we used the number of protein groups, as defined by the Proline software,
288 which include both the unambiguous identification of a single protein (single-protein
289 groups) and groups of indiscernible proteins identified by the same sets of peptides (multi-
290 protein groups) (see “Proteome analysis” in Methods). Interestingly, the proportion of
291 single-protein groups is considerably higher for the reduced database (**Figure 4B**),
292 meaning that protein identifications are less ambiguous.

293 We further characterized ambiguity of protein identifications using the graph’s connected
294 components. Briefly, we first represented peptide-to-protein mappings via bipartite graphs,

295 with peptides and proteins as vertices and with edges featuring peptide to protein
296 membership: this allows an easy picturing of the complex structures generated by shared
297 peptides, as well as their processing by means of graph theory. Then, we calculated
298 connected components (CCs), *i.e.*, the largest subgraphs in which any two vertices are
299 connected to each other by a path and not connected to any other of the vertices in the
300 supergraph. Proteins sharing one or more peptides are thus gathered in the same CC
301 (multi-protein CCs), while unambiguous protein identifications are represented by CCs with
302 a single protein vertex (single-protein CCs) (**Figure 4C**). As such, CCs constitute a
303 peptide-centric strategy to represent ambiguous protein identifications and their shared
304 peptides, not to be confused with the classical protein-centric strategy of protein grouping.
305 It presents two main advantages. First, it provides a non-redundant representation of
306 shared peptides, which can instead be duplicated between different protein groups.
307 Second, it is independent from the different existing strategies of protein inference and
308 protein grouping, making it widely applicable, reproducible and transparent. We observed
309 that, while the total number of obtained CCs is comparable, there is a considerably higher
310 proportion of single-protein CCs in the graph derived from the reduced database search
311 results. After the reduction of the protein group size, this is the second evidence of a
312 decreased ambiguity of protein identifications (**Figure 4D**). Consistently, we also observed
313 a greater proportion of specific peptides – and a correspondingly lower proportion of
314 shared peptides-- from the reduced database search (**Figure 4D**). Within multi-protein
315 CCs, the ratio between the number of protein members and the corresponding number of
316 their encoding genes is also inferior for the reduced database, suggesting that at least part
317 of the solved ambiguity occurred between proteins encoded by different genes (**Figure 4E,**
318 **Supp. Table 6**). As a side note, we additionally observed that searches against reduced
319 databases are associated with inferior ambiguity at the PSM level, although to a less

320 extent. In the peptide identification step, it is common to only consider the best peptide
321 match for each spectrum (*i.e.*, the rank 1 PSM, according to the search engine score) but it
322 can occur that a spectrum matches different peptides equally well (or almost equally). This
323 complicates the analysis and no consensus exists on how to treat these cases.
324 Interestingly, we observed that a smaller proportion of spectra with multiple best matches
325 occur in the reduced database search (**Supp. Fig. 10A**); likewise, fewer best matches are
326 in general found per spectrum (**Supp. Fig. 10B**).

327 At last, we adopted an alternative strategy to enhance proteomics by
328 transcriptomics, which consists in an MS/MS search against the full database, followed by
329 post-hoc filtering of proteins with no expressed transcript and no specific peptide (**Figure**
330 **5A, Supp. Fig. 11**). The driving principle is indeed to remove ambiguous protein
331 identifications not supported by specific peptides or by transcriptomics and thus reduce
332 ambiguity due to shared peptides. Overall, we observed similar results from reduced
333 transcriptome-informed database searches and post-hoc filtering. First, they provide a
334 similar number of spectra and peptide identifications, comparable to that of the full
335 database search (**Figure 5B**); secondly, they yield a similarly higher proportion of single-
336 protein CCs and specific peptides than full database searches (**Figure 5C**), indicating less
337 ambiguous protein identifications. Post-hoc filtering is a transparent and easily
338 interpretable approach and we believe that it is most suitable to studies aiming to enhance
339 protein inference. While a few software tools already exist to generate reduced protein
340 databases, we provide here a specific toolbox of R scripts to perform the aforementioned
341 post-hoc filtering. The toolbox additionally allows very efficient calculation of the CCs,
342 which we have proposed as a means to quantify and compare ambiguity of protein
343 identifications, to visualize CCs of interest and manually inspect them before and after
344 post-hoc filtering.

345

346

347

DISCUSSION

348 In this work, we provide guidance for a mindful use of reduced transcriptome-
349 informed protein databases for MS/MS search. This type of reduced databases stems from
350 the attempt to counter excessive database inflation in proteogenomic studies, when adding
351 variant or novel proteoforms identified from sequencing data. Indeed, increased database
352 size complicates the task of discriminating between correct and incorrect matches. When
353 using TDC-based FDR control, inflated target databases come with an inflated number of
354 decoys and consequently higher probability to get high-scoring decoy matches. This has
355 mainly been thought to reduce sensitivity of identifications in two ways. First, decoy
356 matches may score better than correct target matches and outcompete them in spectrum-
357 peptide assignment (“outcompeting decoys”), so that the number of obtained
358 identifications decreases. Second, decoys may have higher probability to be matched than
359 incorrect targets, which violates the Equal Chance assumption of TDC procedure and
360 provides an overestimated FDR, again decreasing the number of identifications. As the
361 main *raison d’être* of proteogenomics is to maximize the number of identifiable peptides,
362 including variants or non-canonical ones, many efforts were made to avoid loss of
363 sensitivity from excessively large databases, for example by reducing their size. While
364 issues coming from use of excessively large databases have been abundantly discussed,
365 fewer works pointed out that also excessively small databases may be problematic, as
366 they also affect TDC estimations^{31,34,35}. With excessively small databases, TDC provides
367 inaccurate FDR estimates, since these estimates can only be asymptotically
368 accurate^{34,36,37}. Further, with too few (high-scoring) decoys, the probability to match a

369 decoy may be lower than the probability to match an incorrect target, which violates again
370 the Equal Chance assumption of TDC, leading this time to FDR underestimation and to an
371 artifactual increase of identifications. In this work, we showed that the increment of
372 identifications obtained by searching reduced transcriptome-informed protein databases is
373 likely to represent a statistical artifact from employing TDC on excessively small
374 databases. We illustrated how TDC estimates a lower score cutoff for 1% FDR control on
375 the reduced databases compared to the full database search results, causing some invalid
376 PSMs in the full database to be retained as valid additional identifications only in the
377 reduced database. We confirmed this observation at various levels of FDR control (0.5%,
378 1%, 5%) and in two different human-derived samples – a tissue (lung) and a cell line
379 (Jurkat) – with a different level of proteomic complexity and number of spectra. Fewer
380 spectra are available for the Jurkat sample, which, interestingly, also presents a more
381 important difference in score cutoffs between full and reduced database. Indeed, not only
382 reduced database sizes but also a lower number of spectra is believed to affect TDC ability
383 to accurately estimate FDR³⁴. We claim that additional identifications obtained from such
384 reduced databases are at least doubtful and that it is unwise to employ reduced
385 transcriptome-informed protein databases with the aim of increasing the number of
386 identifications. Indeed, the obtained additional identifications have quite low scores and do
387 not stem from removal of out-competing decoys, a known cause of missed identifications
388 in excessively large databases; instead, they rather represent PSMs identical in the two
389 database searches but only accepted in the reduced database due to a lower score cutoff
390 for the same level of FDR control. Most importantly, only a negligible number of additional
391 identifications is generated from the reduced database search when using a method for
392 FDR control known to be stable with respect to database size, such as BH. Indeed, using
393 BH, score cutoffs estimated for the full and reduced database searches, at the same level

394 of FDR control, are almost identical. As is, BH procedure constitutes an interesting
395 alternative to TDC for stable FDR control irrespective of the database size. However, many
396 alternative approaches have been recently developed to cope for the weaknesses of
397 classical TDC^{38–43}. It is important that proteogenomics researchers use them to avoid
398 risking statistical artifacts in their data. In doing so, they will not benefit any longer from the
399 so far hypothesized sensitivity increment, but this seems to be the necessary cost for
400 rigorous control of the FDR.

401 Reduced transcriptome-informed protein databases are nonetheless useful in
402 bottom-up proteomics to reduce ambiguity of protein identifications, which comes from the
403 presence of shared peptides. In particular, we showed that searching these reduced
404 databases yields a higher proportion of specific peptides and unambiguously identified
405 proteins (*i.e.*, single-protein CCs). Furthermore, the higher proportion of specific peptides
406 and correspondingly lower proportion of shared peptides positively affects precision in
407 relative protein quantification. Indeed, in relative protein quantification, where peptide
408 abundances are used as a proxy for the abundance of their parent protein, shared
409 peptides are difficult to handle: since their relative abundance may depend on the
410 contribution of multiple proteins they are frequently discarded. As a downside, this heavily
411 restricts the number of remaining quantifiable proteins, which is reduced to proteins with at
412 least one specific peptide, and the amount of information available to estimate
413 abundances, corresponding to the number of specific peptides only. Therefore, a lower
414 proportion of shared peptides represents more information available for quantification.

415 Finally, we showed that full database searches followed by post-hoc filtering of
416 proteins with no expressed transcript provide proteomic identifications comparable to
417 reduced database searches and similarly reduces ambiguity of protein identifications,
418 while being more transparent and interpretable. We provided an R code to implement such

419 post-hoc filtering strategy. The code allows the user to visualize ambiguous protein
420 identifications and their peptides via bipartite graphs, to prune them according to transcript
421 expression and to manually inspect how this transcriptome-based post-hoc filtering
422 strategy reduces ambiguity. Ambiguous protein identifications are represented and
423 quantified using graph connected components, which constitute here subgraphs of
424 proteins connected by shared peptides. This representation comes with the following
425 advantages: it is transparent, interpretable, non-redundant with respect to shared peptides
426 and independent from the variety of different strategies developed to define protein
427 groups.

428 Results from this work are of interest also beyond proteogenomics. Indeed,
429 database reduction is widely pleaded in proteomics, while little attention is being paid to
430 the limitations of TDC when using excessively small databases. It was proposed for
431 example, to limit database size based on peptide detectability³⁰ and it was more generally
432 claimed that “mass spectrometrists should only search for peptides they care about”⁴⁴. The
433 observed TDC statistical artifacts with excessively small databases is an issue similarly
434 concerning multi-step search strategies implemented by some proteomics search
435 engines⁴⁵ or developed in metaproteomics^{46,47}. Furthermore, the observation that
436 transcriptome information can aid to decrease ambiguity in protein identification is of
437 general relevance in classical proteomics and even more in metaproteomics, which is
438 confronted with an additional source of protein ambiguity, namely the presence of multiple
439 organisms in the same sample.

440

441

METHODS

442

443 **Proteogenomic datasets description**

444 We analysed two samples for which matched transcriptome and proteome were
445 publicly available: a healthy lung tissue and a Jurkat cell line. The lung sample comes from
446 a dataset by Wang *et al.*³³, which includes 29 histologically healthy human tissues and was
447 meant to describe mRNA and protein expression levels across human body. The lung
448 transcriptome dataset was obtained by paired-end RNA sequencing on an Illumina HiSeq
449 2000/2500 system generating 2×100 bases long reads. Its matched proteome dataset was
450 obtained by quantitative label-free LC-MS/MS using an on-line nanoflow liquid
451 chromatography system coupled to a Q Exactive Plus mass spectrometer, operating in
452 data-dependent mode. Sample preparation included peptide fractionation via hSAX
453 (hydrophilic strong anion) chromatography. Transcriptome and proteome raw data were
454 downloaded from the EBI SRA (ArrayExpress accession: E-MTAB-2836; run accession:
455 ERR315346) and the ProteomeExchange (dataset identifier: PXD010154; sample
456 identifier: P013163) repositories, respectively.

457 The Jurkat cell line dataset comes from a study by Sheynkman *et al.*²⁴. The Jurkat
458 transcriptome dataset was obtained by paired-end RNA sequencing on an Illumina HiSeq
459 2000 system generating 2×200 bases long reads. The matched proteome dataset was
460 obtained by quantitative label-free LC-MS/MS using nanoAquity LC system
461 chromatography system coupled to a Velos-Orbitrap mass spectrometer, operating in data-
462 dependent mode. Sample preparation included peptide fractionation via high pH LC
463 separation. Transcriptome and proteome raw data were downloaded from the NCBI's

464 Gene Expression Omnibus (GEO) and the PeptideAtlas repositories with accession
465 GSE45428 and PASS00215, respectively.

466 **Transcriptome analysis**

467 Raw reads were downloaded from public repositories and processed on the Galaxy
468 platform available at <https://usegalaxy.org/>⁴⁸ using common workflows of read
469 preprocessing and alignment for transcript identification (**Supp. Table 3**). First, sequencing
470 adapters and low quality (Phred score < 20) read ends were trimmed off using TrimGalore
471 (https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/) and reads shorter than
472 20 bp after trimming were discarded. Then, preprocessed reads were aligned against the
473 human reference genome (assembly GRCh38) by the splice-aware STAR aligner⁴⁹ in
474 default mode, using the Ensembl reference gene model for splice junctions. Only reads
475 mapped in a proper pair, passing platform quality checks were retained. Reads
476 corresponding to optical or PCR duplicates were removed, as well as non-primary and
477 supplementary alignments. We initially employed two common strategies of transcriptome
478 assembly and quantification: StringTie⁵⁰ and Cufflinks⁵¹. Both programs were run looking
479 for reference transcripts only (no novel transcripts were searched) and they yielded two
480 comparable set of expressed transcripts. Unless otherwise specified, StringTie output was
481 used for downstream analyses.

482

483 **Construction of reduced transcriptome-informed protein databases for**

484 **MS/MS search**

485 For each sample, we built sample-specific protein databases for MS/MS search,
486 containing only those protein sequences for which the corresponding transcript is

487 expressed in the sample. Briefly, we first processed sample-matched transcriptomes as
488 described above and identified the subsets of transcripts expressed at FPKM>1 according
489 to the StringTie or Cufflinks algorithms of transcript assembly and quantification. Then, we
490 filtered the human GRCh38 Ensembl protein database, only keeping those proteins whose
491 corresponding transcript is expressed in the sample. For each sample, we obtained two
492 sample-specific reduced versions of the Ensembl database, based on expressed
493 transcripts from either StringTie or Cufflinks transcript quantification (**Supp. Fig. 1A**).
494 Unless otherwise specified, all downstream analyses were performed using the reduced
495 transcriptome-informed database built according to expressed transcripts identified by
496 StringTie, which is more recent than Cufflinks.

497

498 **Proteome analysis**

499 Raw spectra were downloaded from public repositories and processed
500 automatically using Mascot Distiller software (version 2.7, Matrix Science). Peptide
501 spectrum matches were identified using Mascot search (version 2.6) against two different
502 concatenated target-decoy databases: either the original human GRCh38 Ensembl protein
503 database (release 98, September 2019) or a reduced version of it containing only proteins
504 whose transcript is expressed (as described in the “Construction of reduced transcriptome-
505 based protein databases for MS/MS search”). In both cases an equivalent number of
506 decoy sequences was appended, as well as a custom database of common contaminant
507 sequences (n=500) (and the corresponding number of decoys). Decoy sequences were
508 generated by reversing target sequences with the perl script provided with the Mascot
509 software. The parameters used for Mascot search on the lung and Jurkat samples are
510 reported in **Supp. Table 2**.

511 The Proline software⁵² was used for post-search PSM validation with the following
512 prefilters: *i.* PSMs with score difference < 0.1 were considered of equal score and
513 assigned to the same rank (pretty rank); *ii.* only a single best-scoring PSM is retained per
514 query (single PSM per rank); *iii.* minimum peptide length >= 7 amino acids. Prefiltered
515 PSMs were then filtered at the score cutoff estimated for 1% FDR control. Unless
516 otherwise specified, the score cutoff for FDR control was estimated by target-decoy
517 competition⁵³. No protein inference was performed but for each peptide, all possible
518 protein matches were considered. Protein identifications were reported as protein groups,
519 as defined in the Proline software. Protein groups include both the unambiguous
520 identification of a single protein (single-protein groups) and groups of indiscernible
521 proteins identified by the same sets of peptides (multi-protein groups).

522 Further analyses were performed using the Benjamini-Hochberg procedure for FDR
523 control³², in alternative to TDC (see results section “Transcriptome information aids to
524 reduce ambiguity of protein identifications”). For these analyses, we used PSMs obtained
525 from target-only protein databases, appended with the same database of common
526 contaminant sequences, and searched with the same Mascot parameters as before.

527 **Peptide-protein bipartite graphs and connected components**

528 We represented proteomic identifications using bipartite graphs with two types of
529 nodes -- *i.* identified peptides; *ii.* all their possible proteins of origin -- to more easily
530 analyze and visualize groups of ambiguous protein identifications connected by shared
531 peptides. Indeed, peptide assignments to proteins may be very complex in presence of
532 shared peptides, but are easily represented using bipartite graphs. We then employed
533 graph connected components (CCs), defined as the largest subgraphs in which any two

534 vertices are connected to each other by a path and not connected to any other of the
535 vertices in the supergraph, to quantify and visualize ambiguity of protein identifications.

536 To build bipartite graphs of proteomic identifications, we first generated a tab-
537 separated file containing for each identified peptide all proteins it matches to (one per line),
538 based on the output of PSM validation by the Proline software. We then converted it into
539 an incidence matrix, with proteins along the columns and peptides along the rows, using
540 the `crosstab` function from the GNU `datamash` program
541 (<http://www.gnu.org/software/datamash>). By cross-product of the incidence matrix, we
542 obtained the corresponding adjacency matrix, which describes protein-to-protein
543 connections, based on shared peptides. Finally, we calculated CCs, using the `connComp()`
544 function of the “`graph`” R package on the adjacency matrix. There are two types of CCs: *i.*
545 those containing one single protein (single-protein CCs), with only specific peptides, which
546 constitute unambiguous protein identifications; *ii.* those containing multiple proteins
547 sharing peptides (multi-protein CCs), which represent ambiguous protein identifications.
548 Ambiguous protein identifications can be visually inspected by taking the CC of interest,
549 extracting from the incidence matrix all specific and shared peptides mapping on the CC
550 protein members and plotting peptide-to-protein mappings as bipartite graphs, using the
551 “`igraph`” R package.

552 To decrease the computational cost in case of very large datasets or scarce
553 computational resources, we also developed an alternative strategy of CCs calculation
554 (**Supp. Fig. 12**). First the incidence matrix is reduced by removing all proteins not sharing
555 peptides and all peptides unique to these proteins. Then the corresponding adjacency
556 matrix is generated by cross-product of the incidence matrix and connected components
557 can be more rapidly calculated on this reduced adjacency matrix. In this case, we
558 exclusively obtain multi-protein CCs, since protein identifications with only specific

559 peptides, which correspond to single-protein CCs, were first removed from the incidence
560 matrix. While multi-protein CCs are those of interest when investigating ambiguous protein
561 identifications from shared peptides, single-protein CCs can still be easily retrieved from
562 the original incidence matrix if required.

563 A companion R code is provided. It implements all the above described steps,
564 including: *i.* generating the adjacency matrix; *iii.* calculating connected components; *iii.*
565 visualizing CCs as bipartite graphs.

566

567 **Transcriptome-informed post-hoc filtering**

568 As an alternative to searching a reduced transcriptome-informed database, we
569 tested a transcriptome-informed post-hoc filtering strategy, which works as follows. First,
570 peptide identifications are obtained searching the full reference protein database and
571 validated using the Proline software, as already described in “Proteome analysis” section
572 in Methods. An incidence matrix is generated to encode peptide-to-protein mappings (see
573 “Peptide-protein bipartite graphs and connected components” section). Then, the sample-
574 matched transcriptome is analysed to identify the set of expressed transcripts. Finally, the
575 peptide-to-protein incidence matrix is filtered by removal of proteins with no expressed
576 transcript and no specific peptide. This allows to reduce ambiguity of protein identifications
577 without losing any peptide identification. The one-to-one transcript-to-protein
578 correspondence is guaranteed by the adoption of Ensembl as reference protein database
579 in proteomics and as genome annotation in transcriptomics. The filtered incidence matrix
580 is then converted to an adjacency matrix to calculate CCs, as previously described (see
581 “Peptide-protein bipartite graphs and connected components” section).

582 In this work, we employed the post-hoc filtering strategy on PSMs obtained from
583 searching the target-only full Ensembl protein database followed by Benjamini-Hochberg
584 procedure for FDR control, for comparability with the approach searching target-only
585 reduced transcriptome-informed protein databases, followed by Benjamini-Hochberg
586 procedure for FDR control. Indeed, Benjamini-Hochberg procedure was used in alternative
587 to TDC after searching reduced transcriptome-informed protein databases to obtain
588 accurate FDR control, as explained in the result section (see section “Reduced
589 transcriptome-informed database search does not increase sensitivity if FDR is accurately
590 controlled”). However, in other contexts, post-hoc filtering can be equally well adopted
591 using PSMs from concatenated target-decoy searches followed by TDC-based FDR
592 control. An R code implementing this post-hoc filtering strategy is also provided. In addition
593 to filtering, the code allows visualizing the CCs of interest both before and after post-hoc
594 filtering.

595

596

ABBREVIATIONS

597 BH Benjamini-Hochberg procedure

598 CC connected component

599 DB database

600 FDR false discovery rate

601 MS/MS tandem mass spectrometry

602 PSM peptide-spectrum match

603 TDC target-decoy competition

604

605

DECLARATIONS

606 **Ethics approval and consent to participate**

607 Not applicable

608

609 **Consent for publication**

610 Not applicable

611

612 **Availability of data and materials**

613 All datasets analyzed in this study are publicly available data from the works of
614 Wang et al.³³ and Sheynkman et al.²⁴.

615 The transcriptome datasets were downloaded from the EBI SRA with ArrayExpress
616 accession: E-MTAB-2836 (run accession: ERR315346) and from the NCBI's Gene
617 Expression Omnibus (GEO) with accession GSE45428. The proteome datasets were
618 downloaded from PeptideAtlas with accession PASS00215 and from the
619 ProteomeExchange repository with dataset identifier: PXD010154 (sample identifier:
620 P013163).

621 The code used to perform analyses is available as an R package, under GPL-3
622 license, at <https://github.com/laurafancello/CCs4prot>. Guidelines for installation and usage
623 are documented in the GitHub repository.

624

625 Acknowledgements

626 The authors would like to thank IFB (Institut Français Bioinformatique) infrastructure for
627 providing computational resources and members of the Edyp team for useful comments
628 and feedback.

629

630 Competing interests

631 The authors declare that they have no competing interests.

632

633 REFERENCES

634

- 635 1. Jaffe, J. D., Berg, H. C. & Church, G. M. Proteogenomic mapping as a
636 complementary method to perform genome annotation. *Proteomics* **4**, 59–77 (2004).
- 637 2. Ruggles, K. V *et al.* Methods , Tools and Current Perspectives in Proteogenomics.
638 *Mol Cell Proteomics* **16**,959–981 (2017).
- 639 3. Nesvizhskii, A. I. Proteogenomics: concepts, applications and computational
640 strategies. *Nat. Methods* **11**, 1114–25 (2014).
- 641 4. Zhang, B., Chambers, M. C. & Tabb, D. L. Proteomic parsimony through bipartite
642 graph analysis improves accuracy and transparency. *J. Proteome Res.* **6**, 3549–
643 3557 (2007).
- 644 5. Yang, X. *et al.* DBParser : Web-Based Software for Shotgun Proteomic Data
645 Analyses. *J. Proteome Res.* **3**, 1002–1008 (2004).
- 646 6. Li, Y. F. *et al.* A bayesian approach to protein inference problem in shotgun
647 proteomics. *J. Comput. Biol.* **16**, 1183–1193 (2009).
- 648 7. Nesvizhskii, A. I., Keller, A., Kolker, E. & Aebersold, R. A statistical model for
649 identifying proteins by tandem mass spectrometry. *Anal. Chem.* **75**, 4646–4658
650 (2003).
- 651 8. Serang, O., MacCoss, M. J. & Noble, W. S. Efficient marginalization to compute
652 protein posterior probabilities from shotgun mass spectrometry data. *J. Proteome*
653 *Res.* **9**, 5346–5357 (2010).

- 654 9. Huang, T. & He, Z. A linear programming model for protein inference problem in
655 shotgun proteomics. *Bioinformatics* **28**, 2956–2962 (2012).
- 656 10. Shanmugam, A. K., Yocum, A. K. & Nesvizhskii, A. I. Utility of RNA-seq and GPMDB
657 Protein Observation Frequency for Improving the Sensitivity of Protein Identification
658 by Tandem MS. *J. Proteome Res.* **13**, 4113-9 (2014).
- 659 11. Ramakrishnan, S. R. *et al.* Integrating shotgun proteomics and mRNA expression
660 data to improve protein identification. *Bioinformatics* **25**, 1397–1403 (2009).
- 661 12. Carlyle, B. C. *et al.* Isoform-Level Interpretation of High-Throughput Proteomics Data
662 Enabled by Deep Integration with RNA-seq. *J. Proteome Res.* **17**, 3431–3444
663 (2018).
- 664 13. Bungler, M. K. *et al.* Detection and Validation of Non-synonymous Coding SNPs from
665 Orthogonal Analysis of Shotgun Proteomics Data. *J. Proteome Res.* **6**, 2331–2340
666 (2007).
- 667 14. Alfaro, J. A. *et al.* Detecting protein variants by mass spectrometry: A comprehensive
668 study in cancer cell-lines. *Genome Med.* **9**, 1–12 (2017).
- 669 15. Li, J. *et al.* A Bioinformatics Workflow for Variant Peptide Detection in Shotgun
670 Proteomics. *Mol Cell Proteomics* **10**, 1–11 (2011).
- 671 16. Guerrero-Sanchez, V. M., Maldonado-Alconada, A. M., Sánchez-Lucas, R. & Rey,
672 M.-D. Specific protein database creation from transcriptomics data in nonmodel
673 species: Holm Oak (*Quercus ilex*. L.). *Methods Mol. Biol.* **2139**, 57–68 (2020).
- 674 17. Maringer, K. *et al.* Proteomics informed by transcriptomics for characterising active
675 transposable elements and genome annotation in *Aedes aegypti*. *BMC Genomics*
676 **18**, 101 (2017).
- 677 18. Blakeley, P. & Overton, I. M. Addressing Statistical Biases in Nucleotide-Derived
678 Protein Databases for Proteogenomic Search Strategies. *J. Proteome Res.* **11**,
679 5221–5234 (2012).
- 680 19. Li, H. *et al.* Evaluating the effect of database inflation in proteogenomic search on
681 sensitive and reliable peptide identification. *BMC Genomics* **17**, (2016).
- 682 20. Krug, K., Popic, S., Carpy, A., Taumer, C. & Macek, B. Construction and assessment
683 of individualized proteogenomic databases for large-scale analysis of
684 nonsynonymous single nucleotide variants. *Proteomics* **14**, 2699–2708 (2014).
- 685 21. Park, H. *et al.* Compact variant-rich customized sequence database and a fast and
686 sensitive database search for efficient proteogenomic analyses. *Proteomics* **14**,
687 2742–2749 (2014).
- 688 22. Verbruggen, S. *et al.* Spectral Prediction Features as a Solution for the Search
689 Space Size Problem in Proteogenomics Authors Spectral Prediction Features as a

- 690 Solution for the Search Space Size Problem in Proteogenomics. *Mol Cell*
691 *Proteomics* **20**, 100076 (2021).
- 692 23. Zhang, B. *et al.* Proteogenomic characterization of human colon and rectal cancer.
693 *Nature* **513**, 382–387 (2014).
- 694 24. Sheynkman, G. M., Shortreed, M. R., Frey, B. L. & Smith, L. M. Discovery and Mass
695 Spectrometric Analysis of Novel Splice-junction Peptides Using RNA-seq. *Mol. Cell.*
696 *Proteomics* **12**, 2341–2353 (2013).
- 697 25. Sheynkman, G. M., Shortreed, M. R., Frey, B. L., Scalf, M. & Smith, L. M. Large-
698 scale mass spectrometric detection of variant peptides resulting from non-
699 synonymous nucleotide differences. *J. Proteome Res.* **13**, 228–240 (2014).
- 700 26. Wang, X. *et al.* Protein identification using customized protein sequence databases
701 derived from RNA-Seq data. *J. Proteome Res.* **11**, 1009–1017 (2013).
- 702 27. Wang, X. & Zhang, B. customProDB : an R package to generate customized protein
703 databases from RNA-Seq data for proteomics search. *Bioinformatics* **29**, 3235–3237
704 (2013).
- 705 28. Sheynkman, G. M. *et al.* Using Galaxy-P to leverage RNA-Seq for the discovery of
706 novel protein variations. *BMC Genomics* **15**, 1–9 (2014).
- 707 29. Proffitt, J. M. *et al.* Proteomics in non-human primates : utilizing RNA-Seq data to
708 improve protein identification by mass spectrometry in vervet monkeys. *BMC*
709 *Genomics* **18**, 1–10 (2017).
- 710 30. Shanmugam, A. K., Nesvizhski, A. I., Arbor, A. & Arbor, A. Effective leveraging of
711 targeted search spaces for improving peptide identification in MS/MS based
712 proteomics. *J. Proteome Res.* **14**, 5169–5178 (2015).
- 713 31. Coute, Y., Bruley, C. & Burger, T. Beyond Target – Decoy Competition: Stable
714 Validation of Peptide and Protein Identifications in Mass Spectrometry-Based
715 Discovery Proteomics. *Anal Chem* **92**, 14898-14906 (2020).
- 716 32. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and
717 powerful approach to multiple testing. *J. R. Stat. Soc. Ser. B* **57**, 289–300 (1995).
- 718 33. Wang, D. *et al.* A deep proteome and transcriptome abundance atlas of 29 healthy
719 human tissues. *Mol Syst Biol.* **15**, e8503 (2019).
- 720 34. Gupta, N., Bandeira, N., Keich, U. & Pevzner, P. Target-decoy approach and false
721 discovery rate: when things may go wrong. *J. Am. Soc. Mass Spectrom.* **22**, 1111–20
722 (2011).
- 723 35. Keich, U., Kertesz-farkas, A. & Sta, W. Improved False Discovery Rate Estimation
724 Procedure for Shotgun Proteomics. *J Proteome Res.* **14**, 3148-61(2015).

- 725 36. Levitsky, L. I., Ivanov, M. V, Lobas, A. A. & Gorshkov, M. V. Unbiased False
726 Discovery Rate Estimation for Shotgun Proteomics Based on the Target-Decoy
727 Approach. *J. Proteome Res.* **397**, 6–10 (2017).
- 728 37. Burger, T. Gentle Introduction to the Statistical Foundations of False Discovery Rate
729 in Quantitative Proteomics. *J Proteome Res.* **17**, 12-22 (2018).
- 730 38. Keich, U., Tamura, K. & Noble, W. S. Averaging strategy to reduce variability in
731 target-decoy estimates of false discovery rate. *J. Proteome Res.* **18**, 585–593
732 (2019).
- 733 39. Emery, K., Hasam, S., Noble, W. S. & Keich, U. Multiple Competition-Based FDR
734 Control and Its Application to Peptide Detection. Preprint at
735 <http://arxiv.org/abs/1907.01458v2> (2019).
- 736 40. Yi, X., Gong, F. & Fu, Y. Transfer posterior error probability estimation for peptide
737 identification. *BMC Bioinformatics* **21**, 173 (2020).
- 738 41. Lin, A., Plubell, D. L., Keich, U. & Noble, W. S. Accurately Assigning Peptides to
739 Spectra When Only a Subset of Peptides Are Relevant. *J. Proteome Res.* **20**, 4153–
740 4164 (2021).
- 741 42. Ge, X., Chen, Y. E., Song, D., Mcdermott, M. & Woyshner, K. Clipper : p-value-free
742 FDR control on high-throughput data from two conditions. Preprint at
743 <http://biorxiv.org/content/10.1101/2020.11.19.390773v7.full> (2021).
- 744 43. Etourneau, L., Varoquaux, N. & Burger, T. Unveiling the links between peptide
745 identification and differential analysis FDR controls by means of a practical
746 introduction to knockoff filters. Preprint at
747 <http://biorxiv.org/content/10.1101/2021.08.20.454134v1.full> (2021).
- 748 44. Noble, W. S. Mass spectrometrists should only search for peptides they care about.
749 *Nat. Methods* **12**, 605–608 (2016).
- 750 45. Everett, L. J., Bierl, C. & Master, S. R. Unbiased Statistical Analysis for Multi-Stage
751 Proteomic Search Strategies. *J Proteome Res.* **9**, 700–707 (2010).
- 752 46. Jagtap, P. *et al.* A two-step database search method improves sensitivity in peptide
753 sequence matches for metaproteomics and proteogenomics studies. *Proteomics* **13**,
754 1352–7 (2013).
- 755 47. Cheng, K. *et al.* MetaLab : an automated pipeline for metaproteomic data analysis.
756 *Microbiome.* **5**, 157(2017).
- 757 48. Afgan, E. *et al.* The Galaxy platform for accessible , reproducible and collaborative
758 biomedical analyses : 2018 update. *Nucleic Acids Res.* **46**, W537–W544 (2018).
- 759 49. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21
760 (2013).

- 761 50. Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from
762 RNA-seq reads. *Nat. Biotechnol.* **33**, 290–5 (2015).
- 763 51. Trapnell, C. *et al.* Transcript assembly and quantification by RNA-Seq reveals
764 unannotated transcripts and isoform switching during cell differentiation. *Nat.*
765 *Biotechnol.* **28**, 511–515 (2010).
- 766 52. Bouyssie, D. *et al.* Proline : an efficient and user-friendly software suite for large-
767 scale proteomics. *Bioinformatics.* **36**, 3148–3155 (2020).
- 768 53. Elias, J. E. & Gygi, S. P. Target-decoy search strategy for increased confidence in
769 large-scale protein identifications by mass spectrometry. *Nat Methods.* **4**, 207–214
770 (2007).

771

772

773

774

775

776

777

778

779

780

781

782

783

784

FIGURE LEGENDS

785

786

787 **Figure 1. Comparison of proteomic identifications from the full reference database**

788 **or the transcriptome-informed reduced database searches. A.** Graphical

789 representation of the two MS/MS search strategies we compared. MS/MS searches were

790 performed using Mascot against either the reference human Ensembl protein database

791 (full protein DB) or a subset of the reference database, generated based on transcript

792 expression (reduced protein DB). PSMs were first validated using Proline software with the

793 following prefilters: *i.* PSMs with score difference < 0.1 were considered of equal score and

794 assigned to the same rank (pretty rank); *ii.* only a single best-scoring PSM is retained per

795 query (single PSM per rank); *iii.* minimum peptide length ≥ 7 amino acids ; then they were

796 filtered at the score cutoff estimated by target-decoy competition for 1% FDR control. **B.**

797 Size and overlap of the reference human Ensembl protein database (full protein DB) and

798 the sample-specific reduced transcriptome-informed protein databases (reduced protein

799 DB). **C.** Number of spectra (on the left) or peptides (on the right) exclusively identified in

800 the reduced database (“additional in reduced DB” in blue) or exclusively identified in the

801 full database (“lost in reduced DB” in red) searches. The net difference between additional

802 and lost identifications in the reduced database is also reported on top of each bar (“net”).

803

804

805 **Figure 2. Additional identifications from the reduced database search mainly**

806 **originate from a lower cutoff for 1% FDR control (Jurkat sample). A.** Scatter plot

807 comparing PSMs obtained from the full or reduced database searches. Each data point

808 represents a spectrum: its corresponding PSM score in the full and reduced database

809 searches is reported on the x and y coordinates, respectively. A color code is used to

810 represent the type of match (“target”, “decoy”, or “no match”) for each spectrum in the two
811 searches. Score cutoffs obtained by TDC at 1% FDR are also shown as red and blue lines
812 for the full and reduced database searches, respectively. The up-right insert zooms in on
813 PSMs accepted at 1% FDR only in the reduced database, due to lower score cutoff at 1%
814 FDR (black arrow pointing on the dash circle). **B.** Number of reallocated spectra whose
815 score in the reduced database search is equal to that in the full database or lower. The
816 score from searching the reduced database is never observed to be higher than the score
817 from the full database. **C.** Stripchart reporting the PSM score in the reduced database for
818 spectra undergoing reallocations. Only reallocations to target matches in the reduced
819 database are shown. Reallocations are grouped based on the type of match for the same
820 spectrum in the full and reduced database searches (*<match full DB>_<match reduced*
821 *DB>*). The number and percentage of all spectra in each group is reported on the left. The
822 number of reallocations passing the reduced database cutoff for 1% FDR control is shown
823 in blue: they represent valid reallocations in the reduced database (“nb valid
824 reallocations”). The number of reallocations which would pass the full database cutoff for
825 1% FDR control is shown in red: they represent additional valid identifications exclusively
826 generated by reallocation, independent from the lower score cutoff, and are thus referred
827 to as pure reallocations (“nb valid pure reallocations”). **D.** Bar plot representing the number
828 of spectra (on the left) or the number of spectra identifying additional peptides (on the
829 right) exclusively identified in the reduced database search due to: *i.* lower score cutoff at
830 1% FDR control in the reduced database search compared to the full database; *ii.* pure
831 reallocation. The former are additional identifications from PSMs only passing the cutoff
832 from the reduced database search and which would not be accepted based on the full
833 database cutoff. It includes cases of identical PSMs in both searches (*no reallocation*, in
834 black) and cases of reallocation from decoy (orange), target (gray) or no match (magenta)

835 in the full database search to target matches in the reduced database. Additional
836 identifications from pure reallocation, instead, are those exclusively originated by
837 reallocation, which would also pass the full database cutoff (*i.e.*, independent from the
838 lower score cutoff effect). The Venn diagram on top of the additional peptides graph
839 illustrates the corresponding non-redundant number of additional peptides (*i.e.*, peptides
840 not identified in the full database search) identified from these spectra.

841

842 **Figure 3. Lower cutoff for FDR control in the reduced database to recover valid**
843 **decoys (Jurkat sample). A.** Comparison of valid identifications obtained at 1% FDR from
844 the full database (horizontal red arrow) or reduced database search (vertical blue arrow)
845 and simulation of the valid identifications which would be obtained from the reduced
846 database search if the score cutoff at 1% FDR were equal to that for the full database
847 (dashed red arrow). **B.** Number of valid targets and decoys from the full or reduced
848 database obtained at 1% FDR using the cutoffs estimated by TDC on the respective
849 database search results (first and last rows). The second row instead simulates the
850 number of valid targets and decoys which would be obtained from the reduced database if
851 the estimated cutoff were the same as for the full database. Variations expressed in
852 percentages are shown in gray. The associated nominal FDR level is reported (calculated
853 as $(d+1)/t$, with d and t being the number of valid decoys and targets). **C.** Match in the
854 reduced database search for spectra matching valid targets or valid decoys in the full
855 database. **D.** Score cutoffs obtained by TDC or by BH procedure for FDR control for the
856 full or reduced database searches at various FDR levels (0.5%, 1% and 5%). The variation
857 of score cutoff between full and reduced database searches is reported in percentage.

858

859 **Figure 4. Transcriptome-informed reduced databases yield less ambiguous protein**
860 **identifications. A.** Number of valid identifications obtained from the full (red) or reduced
861 (blue) target-only database searches, followed by BH procedure for 1% FDR control. The
862 number of valid spectra, peptide and protein identifications are reported. Protein groups,
863 as defined by the Proline software, represent here protein identifications and they include:
864 *i.* proteins unambiguously identified by only specific peptides (single-protein protein
865 groups); *ii.* groups of proteins identified by the same set of shared peptides (multi-protein
866 protein groups). **B.** Percentage of single-protein groups. **C.** Bipartite graph representation
867 of peptide-to-protein mappings and usage of graph connected components to visualize
868 and quantify ambiguity of protein identifications. Unambiguous protein identifications are
869 represented by CCs with a single protein vertex (single-protein CCs), while proteins
870 sharing peptides are gathered in the same CC (multi-protein CCs) **D.** Upper panel: total
871 number of connected components. Lower panel: percentage of specific peptides and of
872 single-protein CCs. **E.** Genes encoding proteins from the full and reduced database
873 searches. Upper panel: total number of genes associated to protein matches from the two
874 searches. Lower panel: ratio between the number of protein members in each multi-protein
875 CC and the number of their encoding genes.

876

877 **Figure 5. Transcriptome-informed post-hoc filtering and reduced database search**
878 **strategies similarly reduce protein identification ambiguity. A.** Illustration of the
879 transcriptome-informed post-hoc filtering strategy. First, an MS/MS search was performed
880 against the full canonical protein database. Then, proteins with no corresponding
881 expressed transcript in the sample-matched transcriptome and with no specific peptide
882 (both conditions required) are removed, as well as peptides only mapping to that set of
883 proteins. **B.** Number of valid spectra and peptide identifications obtained from the full or

884 reduced target-only database search (red and blue) or from the post-hoc filtering strategy
885 (orange), after 1% FDR control by BH procedure. **C.** Quantification of protein ambiguity for
886 the full or reduced database search (red and blue) or the post-hoc filtering strategies
887 (orange). Upper panel: total number of obtained CCs. Lower panel: percentage of specific
888 peptides and single-protein CCs.

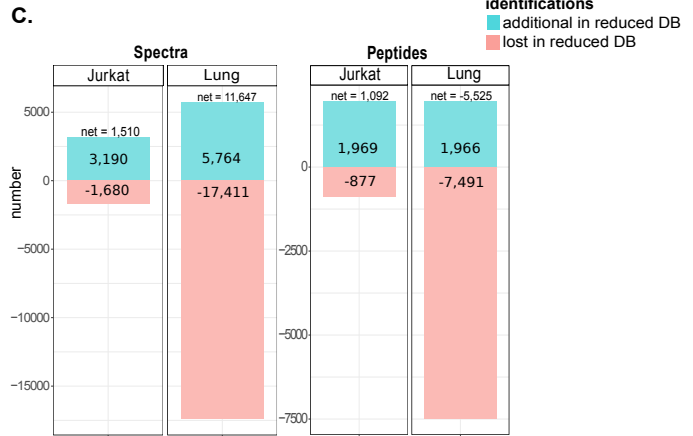
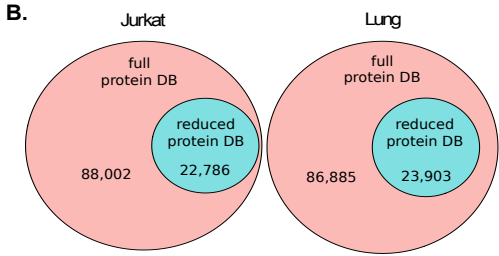
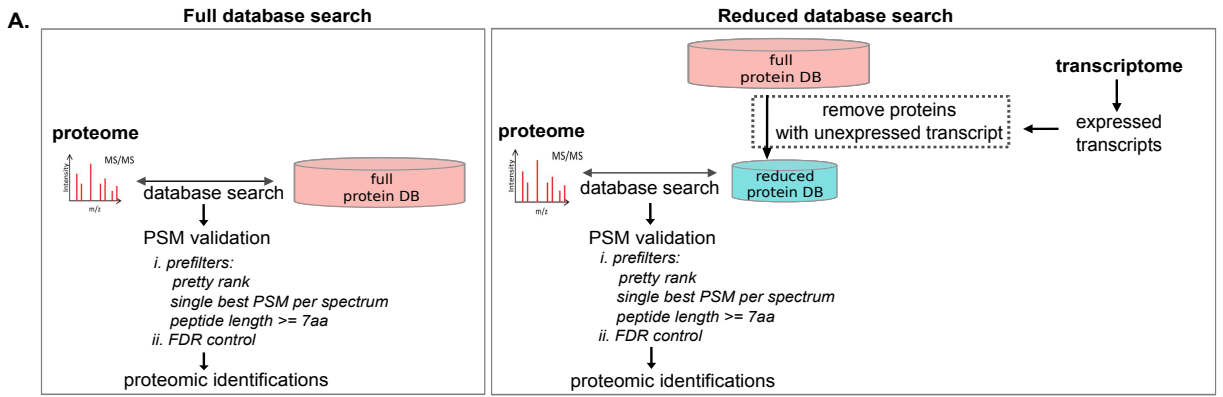
Fig. 1

Fig. 2

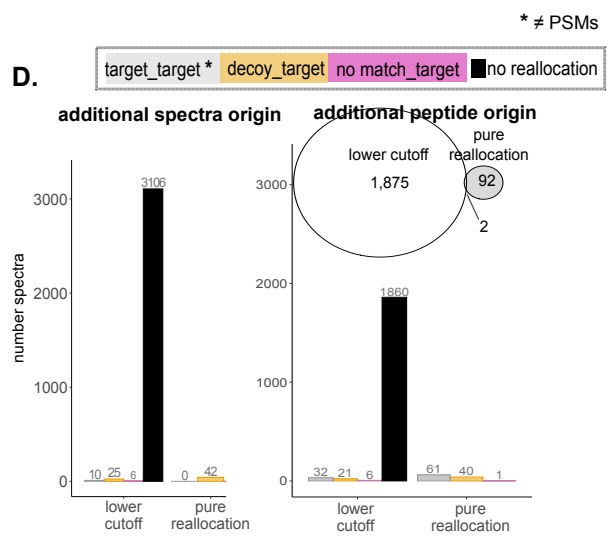
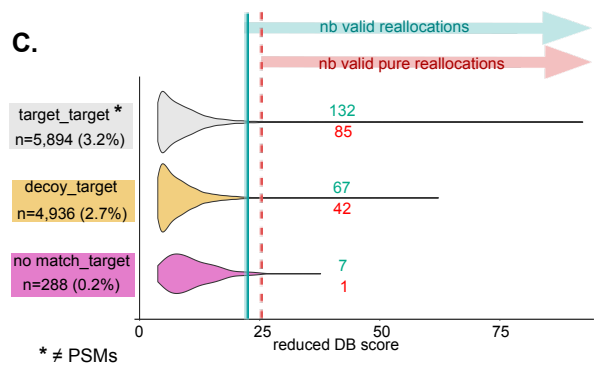
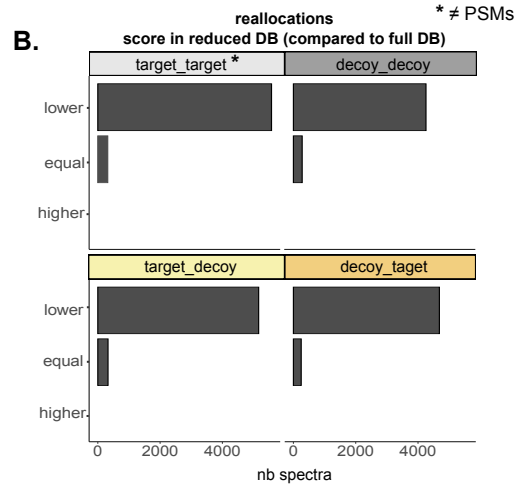
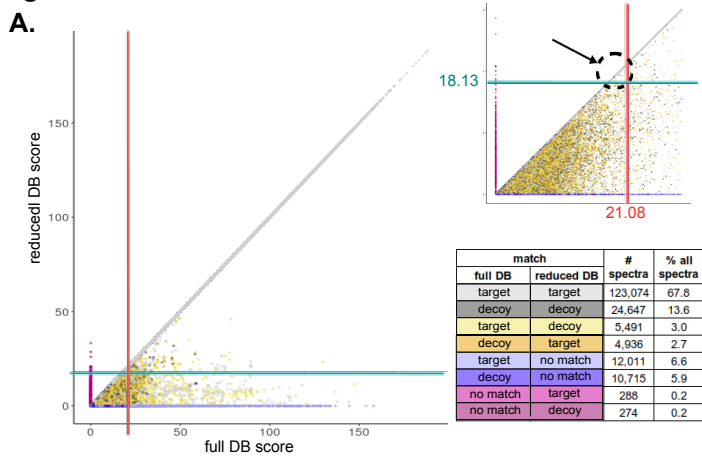
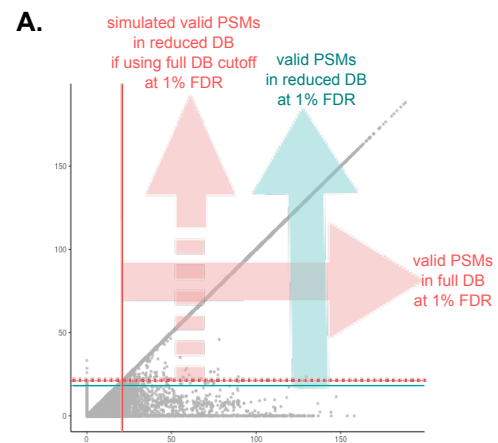


Fig. 3



B.

Database	Score cutoff	# valid targets	# valid decoys	FDR
full	21.08	83,455	833	1
reduced	21.08	81,781	513	0.63
reduced	18.13	84,965	847	1

Changes from full to reduced (21.08): # valid targets -2%, # valid decoys -38.4%

Changes from full to reduced (18.13): # valid targets +3.8%, # valid decoys +40.1%

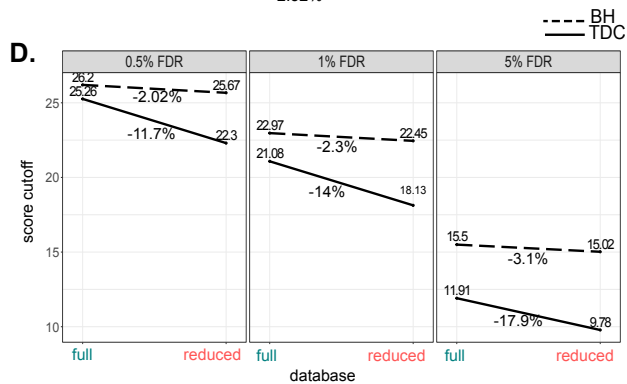
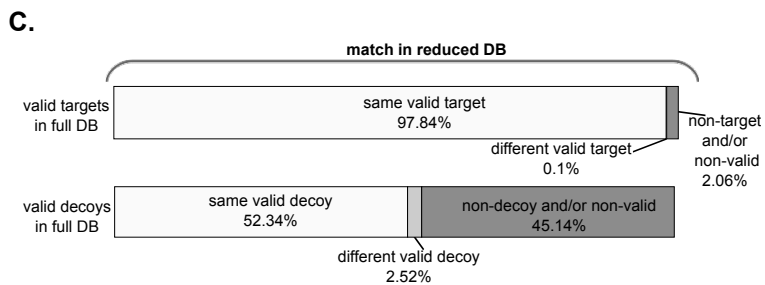


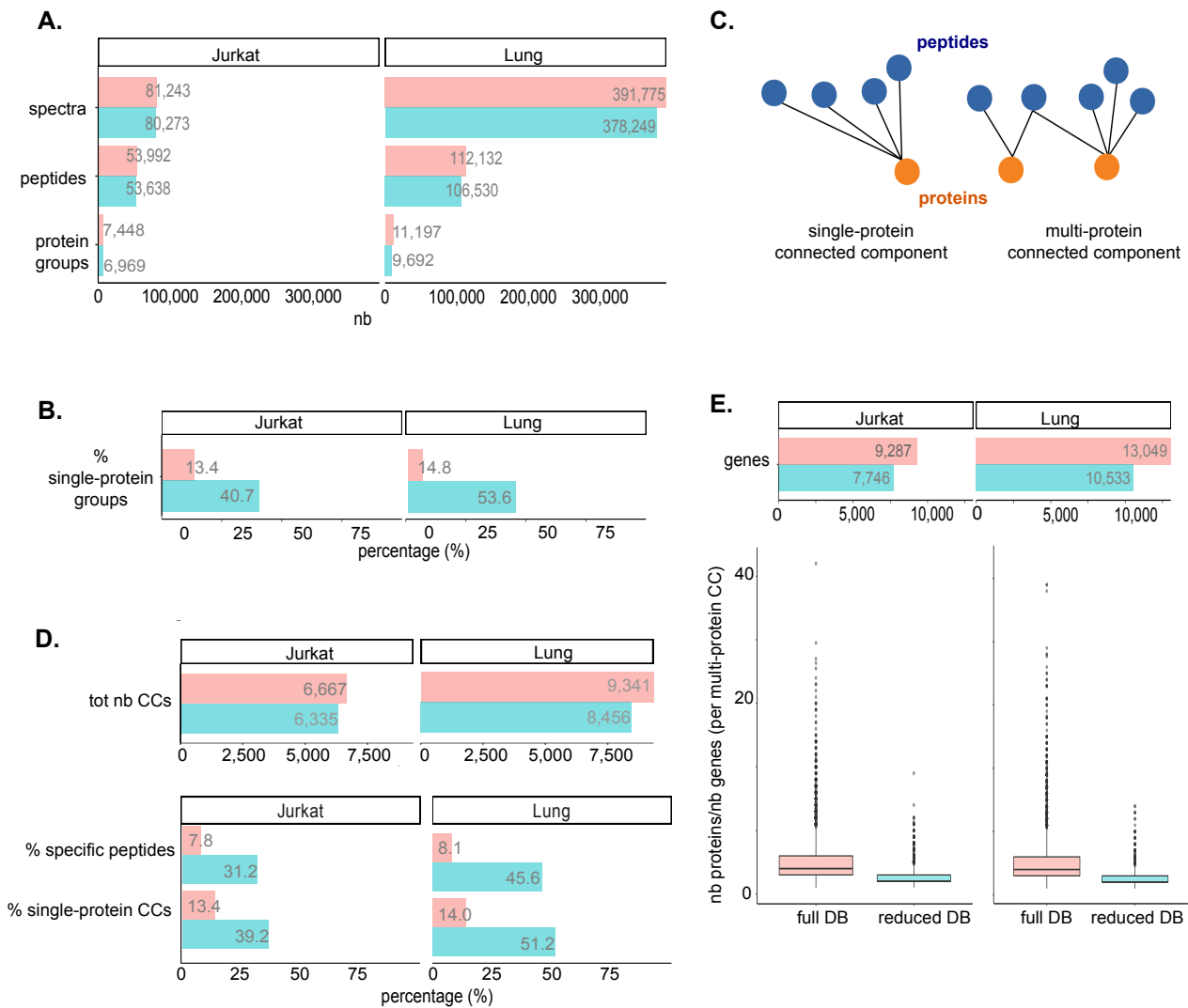
Fig. 4

Fig. 5

