

1 Unifying the analysis of bottom-up proteomics data with CHIMERY5

2

3 **Authors:**

4 Martin Frejno1\$, Michelle T. Berger1\$, Johanna Tüshaus2\$, Alexander Hogrebe1, Florian
5 Seefried1, Michael Graber1, Patroklos Samaras1, Samia Ben Fredj1, Vishal Sukumar1,
6 Layla Eljagh1, Igor Brohnshtein1, Lizi Mamisashvili1, Markus Schneider1, Siegfried
7 Gessulat1, Tobias Schmidt1, Bernhard Kuster2,3, Daniel P. Zolg1 and Mathias Wilhelm3,4*

8

9 \$ Contributed equally

10

11 **Corresponding authors:**

12 *Martin Frejno, martin.frejno@msaid.de

13 *Mathias Wilhelm, mathias.wilhelm@tum.de

14

15 **Affiliations:**

16 ¹MSAID GmbH, Garching b.München, Germany

17 ²Technical University of Munich, Chair for Proteomics and Bioanalytics, Freising, Germany

18 ³Munich Data Science Institute (MDSI), Technical University Munich, Garching, Germany

19 ⁴Technical University of Munich, Computational Mass Spectrometry, Freising, Germany

20

21 **Keywords**

22 CHIMERYs, unifying search algorithm, proteomics, protein identification, protein
23 quantification, proteomics search engine, tandem mass spectra deconvolution, deep
24 learning, DDA, DIA, PRM

25

26 Abbreviations

27	ACN	Acetonitrile
28	AGC	Automatic gain control
29	AWS	Amazon Web Services
30	CAA	Chloroacetamide
31	CID	Collision induced dissociation; synonym for resonance-type CID
32	CV	Coefficient of variation
33	DDA	Data-dependent acquisition
34	DMSO	Dimethyl sulfoxide
35	DIA	Data-independent acquisition
36	FDR	False discovery rate
37	eFDR	Entrapment false discovery rate
38	FA	Formic acid
39	FBS	Fetal bovine serum
40	FFPE	Formalin-fixed paraffin-embedded tissue
41	FTMS	Fourier-transform mass spectrometry
42	FWHM	Full width at half maximum
43	GRU	Gated recurrent unit
44	HCD	Higher-energy collisional dissociation; synonym for beam-type CID
45	ITMS	Ion-trap mass spectrometry
46	I/L isomers	Leucine/Isoleucine isomers
47	MS1	Precursor mass spectrum
48	MS2	Tandem mass spectrum
49	m/z	Mass over charge
50	PBS	Phosphat buffer saline
51	PRM	Parallel reaction monitoring

52	PSM	Peptide spectrum match
53	PD	Proteome Discoverer Software
54	RT	Retention time
55	SILAC	Stable isotope labeling by amino acids in cell culture
56	SPD	Samples per day
57	SVM	Support vector machine
58	Th	Thomson, unit
59	TMT, TMTpro	Tandem mass tags
60	WWA	Wide window Acquisition, synonym for wide window DDA (wwDDA)
61	wwDDA	Wide window DDA, synonym for wide window acquisition (WWA)
62	XIC	Extracted ion chromatogram

63 Abstract

64 Proteomic workflows generate vastly complex peptide mixtures that are analyzed by liquid
65 chromatography-tandem mass spectrometry (LC-MS/MS), creating thousands of spectra,
66 most of which are chimeric and contain fragment ions from more than one peptide. Because
67 of differences in data acquisition strategies such as data-dependent (DDA), data-
68 independent (DIA) or parallel reaction monitoring (PRM), separate software packages
69 employing different analysis concepts are used for peptide identification and quantification,
70 even though the underlying information is principally the same. Here, we introduce
71 CHIMERYYS, a novel, spectrum-centric search algorithm designed for the deconvolution of
72 chimeric spectra that unifies proteomic data analysis. Using accurate predictions of peptide
73 retention time, fragment ion intensities and applying regularized linear regression, it explains
74 as much fragment ion intensity as possible with as few peptides as possible. Together with
75 rigorous false discovery rate control, CHIMERYYS accurately identifies and quantifies multiple
76 peptides per tandem mass spectrum in DDA, DIA and PRM experiments.

77 Main

78 Introduction

79 Mass spectrometry-based bottom-up proteomics is the mainstay technology for high-
80 throughput protein identification and quantification today¹⁻³. The former is achieved by
81 matching theoretical, predicted or library fragment ion mass spectra (MS2) to experimental
82 MS2 spectra, which contain sequence and amino acid modification information on peptide
83 precursor ions, measured in MS1 spectra. Today, MS2 spectra are typically acquired in data-
84 dependent (DDA), data-independent (DIA) or parallel reaction monitoring (PRM) mode.
85 Peptide quantification either uses the peptide ion intensity from MS1 (DDA) or fragment ion
86 intensities from MS2 (DIA, PRM) spectra. A central challenge for data analysis lies in the fact
87 that most MS2 spectra are chimeric, i. e. they contain more than one peptide because LC-
88 MS/MS systems cannot fully separate the vast number of peptides resulting from whole
89 proteome enzymatic digestion.

90 DIA MS2 spectra are usually more complex than DDA MS2 spectra because they are
91 typically acquired with wider isolation windows to maintain low MS cycle times (important for
92 quantification) and hence contain fragment ions from many different precursors⁴. Although
93 DDA and PRM MS2 spectra are typically acquired to minimize co-isolation, they are also
94 chimeric, albeit to a much lesser extent⁵. Because of the way different data acquisition
95 approaches have evolved, the resulting data types are analyzed differently⁶, making it difficult
96 to compare them in an unbiased fashion⁷.

97 DDA data is analyzed in a so-called spectrum-centric fashion⁶. Database search algorithms
98 for DDA data attempt to maximize identifications from chimeric spectra by submitting them
99 multiple times using several precursors detected in the isolation window. Often, fragment
100 ions explained by a given peptide are removed from the spectrum before it is searched
101 again^{5,8}. While often able to call a second or third peptide, this approach will lead to an
102 underutilization of the spectral information when fragments are shared between peptides,
103 resulting in reduced sensitivity. In case fragment ions are not removed for an additional
104 search, there is a danger that the same information is used too often, resulting in reduced
105 specificity. In any case, the central output of DDA search engines is one or multiple peptide-
106 spectrum matches (PSMs) per experimental MS2 spectrum.

107 In contrast, DIA and PRM data analysis follows a so-called peptide-centric approach that
108 asks the question if any of a pre-defined list of peptides are detectable in extracted ion
109 chromatograms (XICs) of their MS1 and/or MS2 spectra^{6,9}. This approach requires the use of
110 peptide spectral libraries, which can be generated from previous experimental data, predicted
111 via machine or deep learning models, constructed directly from the DIA data itself by scoring
112 (deconstructed) MS2 spectra in a DDA-like fashion¹⁰ or a combination of these approaches.
113 Subsequently, the queried peptides are detected and quantified by extracting co-eluting
114 fragment ion chromatograms based on the spectral library.

115 Because of the molecular complexity of proteomic samples and the large quantities of MS2
116 spectra of varying quality that are generated by LC-MS/MS, accurate false discovery rate
117 (FDR) control is an important part of data analysis, particularly in large-scale projects. While
118 FDR control for DDA data is rather mature¹¹⁻¹⁴, it is still a substantial challenge in DIA data,
119 because constructing realistic decoy MS2 spectra and retention times is far from obvious, an
120 issue increasingly realized and addressed by machine learning algorithms for peptide
121 property prediction¹⁵⁻¹⁷.

122 In this work, we introduce a novel spectrum-centric and data acquisition-agnostic approach
123 for the analysis of MS2 spectra, implemented in the search algorithm CHIMERYYS. It
124 deconvolutes any MS2 spectrum, regardless of whether it was acquired by DDA, DIA or
125 PRM, thus unifying the analysis of bottom-up proteomics data. We build upon a concept
126 introduced for the deconvolution of DIA spectra using spectral libraries⁴ and leverage deep
127 learning-based predictions of fragment ion intensities in conjunction with linear algebra for
128 the deconvolution of MS2 spectra. The resulting signal contributions of each peptide
129 identified in each MS2 spectrum can be combined into a quantitative readout. Applying the
130 approach substantially enhances identification rates of PSMs, peptides, and proteins across
131 all sample types in DDA, enables the hands-off processing of PRM data and matches the
132 performance of alternative DIA software while maintaining accurate FDR control throughout.

133 Results

134 Deconvolution of chimeric DDA spectra

135 The core assumption behind CHIMERYYS is that chimeric MS2 spectra are linear
136 combinations of pure spectra from co-isolated precursors. The algorithm is entirely spectrum-
137 centric and employs non-negative L1-regularized regression via the LASSO¹⁸ to explain as
138 much experimental intensity as possible with as few peptide precursors as possible (Figure
139 1A). It uses highly accurate predictions of fragment ion intensities and retention times for
140 target and decoy peptides instead of spectral libraries.

141 Briefly, predicted MS2 spectra from precursors with predicted retention times that fall within a
142 data-dependent retention time window and precursor isotope envelopes that (partially)
143 overlap with the isolation window are compared to experimental MS2 spectra based on a
144 multitude of fragment ion intensity-free and -dependent scores for each PSM (Online
145 Methods). Next, spurious PSMs are removed based on some of these scores. For example,
146 PSMs are required to have at least three matched fragment ions, one of which must be the
147 base peak (most abundant peak of the prediction) and another one of which must be among
148 the top three most intense peaks of the predicted spectrum. PSMs passing these criteria are
149 used for deconvolution, where they compete for experimental fragment ion intensity in one
150 concerted step; an approach fundamentally different from the classic subtraction methods
151 (Figure 1A). PSMs with enough contribution to the experimental spectrum as measured by
152 CHIMERYYS coefficients and that pass additional score filters are handed to mokapot¹³ for
153 PSM-level FDR control, specifically allowing for multiple PSMs per spectrum, similar to
154 DIAMeter¹⁹.

155 We validated this FDR estimation on data with varying chimericity by systematically
156 increasing the isolation window widths of 1-hour HeLa single-shot measurements from 1.4 up
157 to 20.4 Th using entrapment experiments (Online Methods). Figure 1B shows that
158 CHIMERYYS' peptide-level q-values correspond to empirical q-values calculated based on
159 entrapment identifications with the classic eFDR approach, independent of isolation window
160 width.

161 Figure 1C displays the confident identification of six peptides with relative contributions to the
162 experimental total ion current ranging from 4% to 54% in a mirror spectrum. Notably, the
163 experimental intensities for the y1, y1-NH₃ and y1-H₂O ions that are shared between the five
164 peptides (C-terminal lysine) align well with the sum of predicted intensities of the
165 corresponding peptides, scaled by their respective CHIMERYYS coefficient, which can be
166 interpreted as the interference-corrected total ion current of a peptide in an MS2 spectrum
167 (Online Methods). This exemplifies how the algorithm identifies multiple peptides in chimeric

168 spectra while distributing intensities of shared fragment ions. Peptides identified by
169 CHIMERYYS recapitulate the expected quantitative ratios in a multi-organism-mixture
170 experiment (Figure 1D). This renders CHIMERYYS suitable for approaches like wide-window
171 DDA acquisition (also termed WWA or wwDDA)^{20,21} and the direct analysis of DIA data.

172 To assess the performance of the algorithm on DDA data, we analyzed a 2-hour HeLa cell
173 lysate digest with 1.3 Th MS2 isolation windows (Online Methods). CHIMERYYS identified
174 238,795 PSMs at 1% PSM FDR with an overall identification rate of >85% (Supplementary
175 Figure 1A). More than two thirds of identified MS2 spectra contained more than one
176 precursor (Supplementary Figure 1B) confirming previous observations⁵. Fragment ions
177 shared between different peptides were detected across the full MS2 m/z range with an
178 expected higher frequency <200 m/z (Supplementary Figure 2), rendering current
179 approaches for handling chimeric spectra error prone. Comparing these results to seven
180 academic and commercial DDA search engines (Figure 1E) revealed that CHIMERYYS
181 identifies many additional peptides (Supplementary Figure 3). Most of these additional
182 identifications stem from low abundant peptides (Supplementary Figure 4A) with fewer
183 matched fragment ions (Supplementary Figure 4B) that were identified based on intensity-
184 dependent scores such as the normalized spectral contrast angle (Supplementary Figure
185 4C). This resulted in a markedly higher number of peptides per protein group in CHIMERYYS
186 compared to Sequest HT (Supplementary Figure 4D). It is worth noting that some of these
187 search engines do not control FDR at the same level, which has a substantial influence on
188 such comparisons (Supplementary Figure 4E-F). Controlling FDR at a 'lower' level and
189 counting identifications at a 'higher' level (e.g. counting peptides at PSM FDR) will
190 overestimate the number of identifications. Identifications need to be reported at the same
191 level at which FDR is controlled.

192 The gains observed for HeLa digests relative to Sequest HT were corroborated using more
193 difficult biological samples at the protein group level (urine: +21%; CSF: +17%; plasma:
194 +10%; FFPE material: +37%; secretomes: between +33% and +71%, *Arabidopsis thaliana*:
195 +13%; *Haliobacterium*: +20%) (Supplementary Figure 5A-F). This data highlights that
196 CHIMERYYS substantially increases the analysis depth of DDA data without changing data
197 acquisition.

198 Revisiting legacy data using CHIMERYYS

199 We conducted a retrospective study of HeLa single-shot analyses spanning many years and
200 Orbitrap instrument generations. Despite many differences that impair a truly fair comparison,
201 a clear trend was observed in that the higher the speed and sensitivity of the instrument, the
202 higher the advantage of CHIMERYYS over Sequest HT (Figure 2A, Supplementary Figure
203 6A).

204 Next, we investigated low-resolution ion trap data (ITMS) comparing CHIMERYYS to Sequest
205 HT on unprocessed spectra and on spectra filtered for containing only the top 15 most
206 abundant fragments per 100 Th window (Figure 2B). In contrast to high-resolution Orbitrap
207 data, we observed a notable improvement by removing low-abundance peaks in ITMS
208 spectra. Specifically, CHIMERYYS identified 74% more PSMs, 35% more peptides, and 30%
209 more protein groups compared to Sequest HT on unprocessed spectra were, while it
210 identified 94% more PSMs, 47% more peptides, and 37% more protein groups on spectra
211 preprocessed with a top 15 by 100 Th filter. Both examples show that substantially more
212 information can be extracted from legacy data by harnessing the information contained in
213 chimeric spectra.

214 [Optimizing data acquisition with deconvolution in mind](#)

215 We next assessed to what extent CHIMERYYS' capability to deconvolute highly complex
216 spectra can be used to optimize data acquisition. First, we evaluated LC gradients with the
217 goal to increase sample throughput per day (SPD; Online Methods). Figure 2C shows that
218 CHIMERYYS identified a similar number of peptides and proteins in 30 min (48 SPD) for which
219 Sequest HT needed 120 min (12 SPD) of the same HeLa digest, increasing throughput by a
220 factor of four.

221 Next, we explored a possible increase in identification efficiency by intentionally widening the
222 isolation windows in DDA (between 1.4 Th to 20.4 Th; Supplementary Figure 6B-D). The
223 analysis revealed that the number of identified PSMs increased with wider isolation windows
224 (Figure 2E) and began to plateau at > 8 m/z. This is likely due to the AGC limit, which –
225 together with the dynamic range of MS2 spectra – limits the number of peptides in chimeric
226 spectra with a sufficient number of detectable fragment ions. The number of peptide (and
227 protein) identifications reached its maximum at a window size of 3.4 Th for this specific
228 dataset and substantially decreased for larger isolation windows, likely due to the fact that
229 more and more of the limited number of PSMs were from the same, high-abundant peptides.
230 Such approaches have also gained popularity in single-cell proteomics (SCP), where
231 extended injection times enhance sensitivity but result in fewer MS2 scans. CHIMERYYS can
232 counteract this effect and was already applied to extract more PSMs and peptides from these
233 intentionally chimeric scans in SCP data²¹. The strong gains at the PSM, peptide and protein
234 level are noteworthy as identifying nearly 8,000 proteins in a single 120 min DDA analysis
235 has rarely been achieved before and implies that chromatographic pre-fractionation of
236 samples may no longer be necessary.

237 [Deconvolution of chimeric DIA spectra](#)

238 CHIMERYYS deconvolutes DIA spectra in the same way as described for DDA spectra above.
239 The only difference is that DIA spectra are usually more chimeric. Exemplified by a high-load

240 LFQbench-type multi-organism mixture dataset²², CHIMERYYS identified an average of
241 529,993 PSMs per raw file at 1% run-specific PSM FDR, mapping to 66,888 unique peptide
242 groups and 7,331 unique protein groups at 1% global peptide group and protein FDR,
243 respectively, with an overall identification rate of >60% (Supplementary Figure 7A). More
244 than 82% of identified MS2 spectra contained more than one precursor (Supplementary
245 Figure 7B) and shared fragment ions were more frequent, emphasizing the need for
246 spectrum deconvolution that assigns shared fragment ions *pro rata* to the contributing
247 peptides (Supplementary Figure 7C-J).

248 Comparison to other DIA search engines

249 We compared results obtained with CHIMERYYS on DIA data acquired on an Orbitrap QE HF-
250 X from the LFQbench-type dataset to the library-free workflows implemented in the popular
251 software tools DIA-NN²³ and Spectronaut²⁴ using entrapment experiments to validate FDR
252 control in the run-specific context²⁵ (see Online Methods for context definitions and search
253 parameters). The results show that CHIMERYYS' self-reported q-values correspond to the
254 empirical q-values calculated based on entrapment identifications (Supplementary Figure
255 8A). DIA-NN and Spectronaut appeared to underestimate FDR based on all three or the
256 peptide and concatenated entrapment approaches, respectively (Supplementary Figure 8B-
257 C). Recently proposed more stringent settings for Spectronaut²⁶ had little if any effect on this
258 issue (Supplementary Figure 8D). Similar observations were made when analyzing the
259 TimsTOF Pro data of the LFQbench-type dataset using Spectronaut (Supplementary Figure
260 8E). All analyses below used the peptide eFDR approach (eFDR from here onwards).
261 Filtering on eFDR in addition to the algorithm-dependent self-reported FDR did not change
262 the overall number of identifications for CHIMERYYS. Reductions to a level comparable to
263 CHIMERYYS' results were observed for DIA-NN and Spectronaut. Data completeness for
264 CHIMERYYS did not change when requiring precursors to be quantified in two out of three
265 replicate experiments. Reductions to a level similar to CHIMERYYS' results were observed for
266 Spectronaut and to a level below CHIMERYYS for DIA-NN (Figure 3A). At full data
267 completeness, CHIMERYYS and Spectronaut substantially outperformed DIA-NN on the
268 LFQbench-type data (Supplementary Figure 9A).

269 As one might expect, precursors filtered out based on eFDR have lower MS2 intensities
270 (Figure 3B) and fewer fragment ions (Figure 3C). However, the extent to which this is
271 observed differs substantially between the three tools, with CHIMERYYS considering far more
272 fragment ions for quantification than the other two and being more rigorous in the inclusion of
273 fragment ions with very low intensity when using the corresponding default settings. The
274 latter is illustrated in Figure 3D, in which the top panel shows fragment ion chromatograms
275 for a peptide confidently identified by all three search engines, and the bottom panel shows

276 fragment ion chromatograms for a peptide identified only by Spectronaut and for which no
277 evidence of co-elution of fragment ions can be observed (see also Supplementary
278 Information). Further investigations regarding the number and intensity of fragment ions, as
279 well as the corresponding raw data (Supplementary Figure 9B-D) suggest that precursors
280 with less than three quantifiable fragment ions with an intensity exceeding 1 or those with
281 near-zero (or zero) intensity should be removed; either categorically or by applying stringent
282 FDR control, which has a very similar effect (Supplementary Figure 9E). The latter brings all
283 three software tools to a comparable level of overall identifications.

284 [Accurate peptide quantification from chimeric PRM and DIA spectra](#)

285 One of CHIMERY'S distinguishing concepts is its spectrum-centric processing of chimeric
286 spectra. Apart from peptide identification, it also derives spectrum-centric quantitative
287 information in the form of CHIMERY'S coefficients, which can be interpreted as the
288 interference-corrected total ion current for a given peptide in this MS2 spectrum (Online
289 Methods). If none of the matched fragments for a peptide are shared with another peptide
290 and the predicted MS2 spectrum matches perfectly to the experimental one, the coefficient is
291 the sum of all matched fragment ions in the experimental MS2 spectrum. Hence, tracing the
292 coefficient along retention time generates a pseudo-extracted-ion-chromatogram (XICs) that
293 can be used to perform (relative) quantification of peptides based on their MS2 signal in PRM
294 and DIA data. This is different from standard approaches that create XICs for (a subset of)
295 fragment ions of a given peptide, which need to remove interfered fragment ions from
296 quantification to maintain high precision and accuracy (Figure 4A). To assess the
297 performance of our concept, we performed a simple PRM assay, focusing on 52 peptides
298 from 18 human proteins spanning five orders of magnitude of cellular abundance (Online
299 Methods). Both CHIMERY'S and Skyline recovered 47 out of 52 peptides from the targeted
300 inclusion list and CHIMERY'S' automatically-generated MS2-based quantification was in
301 excellent agreement ($R=0.99$) with the manually curated values obtained from Skyline
302 (Figure 4B). Without any additional effort, CHIMERY'S identified and quantified 1,400 further
303 peptides that were not designed to be in the assay but that happened to be co-isolated along
304 with the targeted peptides (Supplementary Figure 10A-C). CHIMERY'S effectively automates
305 the processing of PRM data because it removes the manual curation steps often required in
306 Skyline. These include dealing with shared fragment ions and co-isolated peptides (both
307 used in CHIMERY'S but removed in Skyline).

308 Next, we compared the MS2-level quantitative precision and accuracy of CHIMERY'S to DIA-
309 NN and Spectronaut on the LFQbench-type dataset²². To avoid differences in quantification
310 due to different methods for determining peak integration borders, we compared the three
311 algorithms based on their implementation of peak apex quantification. When filtering the data

312 using eFDR as discussed above, the median quantitative precision of precursors (based on
313 coefficient of variation, CV) was 26.9%, 29.9% and 29.1% for CHIMERYs, DIA-NN and
314 Spectronaut, respectively (Figure 4C, bottom panel).

315 Similarly, analyses of precursor-level ratio distributions (Figure 4D) as a measure of
316 quantitative accuracy for the three different search engines at eFDR were comparable (mean
317 \log_2 -ratios +/- standard deviation for *Escherichia coli*, *Homo sapiens* and *Saccharomyces*
318 *cerevisiae* of -1.90 +/- 0.25, -0.03 +/- 0.25 and 1.00 +/- 0.29 for CHIMERYs, -1.83 +/- 0.32, -
319 0.04 +/- 0.23 and 0.99 +/- 0.28 for DIA-NN and -1.83 +/- 0.35, -0.04 +/- 0.31 and 1.00 +/- 0.37
320 for Spectronaut, respectively). The above analysis demonstrates that CHIMERYs' spectrum-
321 centric way of quantifying peptide precursors matches the performance of Skyline on PRM
322 data as a gold standard in the field and extends to full-scale DIA data. It also highlights the
323 potential of CHIMERYs for scaling PRM assays to very large numbers of peptides without
324 the need for manual intervention.

325 **Head-to-head comparison of DDA and DIA data, facilitated by CHIMERYs**

326 We showed that CHIMERYs can analyze DDA and DIA data using the same concepts for
327 the deconvolution of chimeric spectra, which enables directly comparing the two acquisition
328 methods on data acquired from the same sample, without the need to process the data with
329 different software packages. As one would expect, it identified more than twice as many
330 PSMs from DIA (8 m/z isolation windows) compared to DDA (1.3 m/z isolation windows) data
331 acquired on an Orbitrap QE HF-X (LFQbench-type dataset; Supplementary Figure 11A).
332 However, DDA identified 52% more peptides (and 30.3% more protein groups) compared to
333 DIA (Figure 5A, Supplementary Figure 11B). Likely, this is due to the interplay between the
334 AGC limit and the dynamic range in MS² spectra, which we already observed for WWA data
335 (see section above). In contrast, relative quantitative data completeness was higher for DIA
336 than for DDA data when filtering for peptides that met 1% FDR in the global, but not
337 necessarily the run-specific context and enabling 'match between runs' for DDA using the
338 Minora Feature Detector in Proteome Discoverer²⁷ (78% versus 55.4% of peptides quantified
339 in two out of three replicates per condition in DIA and DDA, respectively, Figure 5A). This
340 resulted in very similar numbers of peptides being quantified in two out of three replicates
341 (56,322 and 52,161) for DDA and DIA, respectively.

342 Perhaps the more interesting comparison is that of DDA vs DIA using the same isolation
343 window (here 2 m/z). This has recently become possible because modern, fast scanning
344 instruments blur the border between DDA and DIA²⁸. Interestingly, both 14 min (~100 SPD)
345 and 30 min (48 SPD) gradients on an Orbitrap Astral²⁹ yielded similar numbers of PSM,
346 peptide and protein group identifications for DDA and DIA (Figure 5A, Supplementary Figure
347 11A-B). The small differences in favor of DIA are likely due to the higher scan rate of the

348 Orbitrap Astral in DIA mode. Again, relative quantitative data completeness was much better
349 for DIA than for DDA (97.9% and 98.7% of peptides quantified in two out of three replicates
350 vs 56.3% and 61.7% for the 14 min and 30 min gradients, respectively; Figure 5A). This data
351 suggests that DIA and MS2-based quantification should be preferred over DDA and MS1-
352 based quantification when performing label-free single-shot measurements on fast scanning
353 instruments. Comparing CV distributions of peptides detected by DDA and DIA in the three
354 datasets revealed that DDA was slightly more precise on the Lfqbench-type dataset, while
355 DIA was slightly more precise on the 30 min Orbitrap Astral dataset (Figure 5B). Quantitative
356 accuracy appeared to be generally better for DIA on the Lfqbench-type dataset (Figure 5C).
357 However, closer inspection suggests that this is due to a problem with the samples rather
358 than with MS1-based quantification *per se*, since the accuracy of MS1- and MS2-based
359 quantification of the DIA data is comparable (Supplementary Figure 11C). In fact,
360 CHIMERYs' MS2-based quantification was highly correlated ($R=0.88$) to the MS1-based
361 quantification implemented in Proteome Discoverer on the same raw data (Figure 5D),
362 suggesting that the two quantification methods could be combined in the future in
363 CHIMERYs.

364 Discussion

365 In many ways CHIMERYs returns to the very old concept of analyzing tandem mass spectra,
366 one at a time. At least for the task of peptide identification, this so-called spectrum-centric
367 approach places the core analytical evidence acquired by the mass spectrometer at the
368 center of all data analysis. This comes with a number of important advantages. First, any
369 proteomic data type (DDA, DIA, PRM) can be treated the same and CHIMERYs is the first
370 software implementation that stringently follows this unifying philosophy. Second, there is no
371 principle difference between identifying a single or multiple peptides from the same MS2
372 spectrum and skilled scientists have done so since the early days of proteomics. The added
373 sophistication is that artificial intelligence can predict the tandem mass spectrum of any
374 peptide with outstanding accuracy so that it is possible to deconvolute even highly chimeric
375 spectra by maximizing the explained intensity in an MS2 spectrum using a minimal set of
376 peptides to do so. Third, statistical methods for PSM-level FDR control are conceptually well
377 worked out and have reached a very high level of practical refinement, again including the
378 use of artificial intelligence that can predict the tandem mass spectrum of any target or decoy
379 peptide with the same accuracy, ensuring fair competition between targets and decoys.
380 Fourth, the plausibility of an identification can be further assessed (albeit not automatically)
381 beyond statistics by visual inspection in the context of the full MS2 spectrum and e. g. looking
382 out for fragment ions that were not part of the deep learning model and have thus not yet
383 been used for identification. A current limitation of CHIMERYs in this context is that peptides

384 carrying modifications that are not yet covered by the underlying deep learning model escape
385 detection. It can be anticipated that this limitation will diminish over time as deep learning
386 models start to emerge that are capable of generalizing to modifications or fragmentation
387 methods that they have not yet been trained for³⁰.

388 Akin to other software tools, CHIMERY5 also uses the information contained in the MS2
389 spectrum for peptide quantification. However, unlike all other DIA software, it does not set a
390 fixed number of fragment ions to consider and instead always uses all the fragment ions that
391 have led to an identification in a given MS2 spectrum, but in relative proportion to how much
392 they contributed to the actual signal in the MS2 spectrum (important for the frequent case of
393 fragment ions that are shared between peptide candidates). CHIMERY5 uses the sum of
394 these fragment ion intensities rather than the individual fragment intensities to find the apex
395 of a chromatographic peak. This makes the overall quantification more robust against weak
396 signals and spurious detections as encountered in e. g. single cell proteomics data. The
397 results indicate that quantitative precision and accuracy closely match that of PRM data,
398 which is often considered to be the gold standard for peptide quantification. In this context it
399 is interesting to note that CHIMERY5 also “en passant” automates the analysis of PRM
400 experiments.

401 We consciously decided to rate data quality over quantity such that reported peptide
402 identification and quantification results are rather conservative and other software tools may
403 sometimes seemingly outperform CHIMERY5 (see also Supplementary Information).
404 However, when applying rigorous and consistent criteria for peptide detection and
405 quantification, these differences diminish. A perhaps unexpected finding in this regard is that
406 DIA data is often not nearly as complete as default processing parameters of DIA search
407 engines report. Again, and not surprisingly, this is particularly true for low abundant samples
408 or low abundant peptides within a sample. The reasons for this could be manifold and
409 investigating them comprehensively goes beyond the scope of the present study. However, it
410 is worth mentioning that the most recent generation of mass spectrometers have driven
411 sensitivity to the point of single ion detection. As a result, MS2 spectra have at least some
412 low level of signal at nearly every m/z. Many of these may not even stem from peptides but
413 will create a situation in which “something” can be easily found everywhere and all the time,
414 leading to data completeness that bears little if any actual justification. In addition, the
415 increasing volume and density of MS-based proteomic data keeps challenging the scalability
416 of the assumptions underlying data processing tools. Reassuringly, the community of
417 proteomics software developers and users are increasingly aware of these recurring
418 challenges, as it is in everybody’s best interest to ensure that software tools can be trusted
419 and used at face value. CHIMERY5 makes a valuable contribution in this context and a
420 particularly exciting prospect is that the latest LC-MS/MS hardware along with the latest

421 software solutions will soon overcome the historically grown divide in the field between DDA
422 and DIA.

423

424 Online Methods

425 Deep learning for peptide property prediction

426 Data preprocessing – fragmentation prediction

427 Publicly available data from various PRIDE identifiers were used as a training data
428 foundation, notably the ProtomeTool project³¹. RAW data were downloaded and – where
429 available – MaxQuant⁸ search files were utilized. PSMs from MaxQuant’s msms.txt were
430 merged with the unprocessed scans extracted from RAW files with Thermo Fisher’s Raw-
431 FileReader (<http://planetorbitrap.com/rawfilereader>). Data were filtered using various quality
432 criteria, e.g. Andromeda score. The top 3 ranked PSMs by Andromeda score per sequence,
433 charge, collision energy combination were selected across all files. For the PSMs, b- and y-
434 ions, as well as several neutral losses (e.g. water, ammonia and carbon monoxide,
435 phosphoric acid losses) were annotated for charge states 1-3. Amino acid tokens in the
436 peptide sequence, precursor charge, and fragmentation type were one-hot encoded, and
437 collision energy was normalized to [0, 1]. The data contains 25M PSMs and was split into
438 70% training, 20% test and 10% validation sets, ensuring that peptide sequences are not
439 shared across splits.

440 Data preprocessing – retention time prediction

441 The same unprocessed data was transformed similarly for retention time training with the
442 following differences: in addition to a stringent Andromeda score threshold, only the top-
443 ranked PSMs by Andromeda score per sequence and file, and only the 5 top-ranked PSMs
444 per sequence across all files were retained. Retention time was z-score normalized per file to
445 account for differing run lengths. The resulting dataset contains 5.2M PSMs split into 70%
446 training, 20% validation, 10% test sets, ensuring that peptide sequences are not shared
447 across splits.

448 Model architecture

449 INFERYS 3.0 models for fragmentation and retention time share a similar architecture.
450 Sequences are processed with a PositionalEmbedding³², a Transformer block and a GRU³³
451 layer (gated recurrent unit). For fragmentation prediction, additional meta information
452 (precursor charge, collision energy and fragmentation type) are injected via a custom
453 TransformerMixin layer to the sequence embedding outputs of the Transformer block (before
454 the GRU layer). The TransformerMixin embeds one input parameter to the dimensionality of
455 a given Transformer output embedding and applies the product of the two to another
456 Transformer layer. The final Transformer embedding is projected to the task-dependent

457 output dimension (e.g. 1 output dimension for retention time). Models are built and trained
458 with tensorflow 2.11.1.

459 **Model training**

460 The same training procedure was used for the fragmentation and retention time model.
461 Models are trained with the Adam optimizer³⁴ on the training split with early stopping,
462 evaluating the validation split with patience 8 and a learning rate decay with a factor of 0.2
463 after 4 epochs without reduced validation loss for up to 200 epochs or till convergence. The
464 retention time model optimizes mean-absolute-error and the fragmentation model optimizes
465 normalized spectral contrast distance^{15,35}. Model hyperparameters, such as learning rate,
466 batch size, dropout, embedding dimension, positional scaling, number of attention heads,
467 and intermediate dimensions are optimized via Hyperband³⁶. For the fragmentation model,
468 we also optimized the order of meta parameter TransformerMixin layers via Hyperband.

469 **Model capabilities**

470 The INFERYS 3.0 fragmentation model predicts a set of fragment ions consisting of b-ions,
471 y-ions, water-loss ions ammonia-loss ions, carbon monoxide-loss ions and phosphoric acid-
472 loss ions in charge states 1 to 3. The exact set of ions was predetermined by selecting ions
473 that explain the most experimental intensity in the training data set, hence some ions never
474 or rarely observed in the training data were excluded (e.g. triply-charged y1). The
475 fragmentation model is compatible with peptides of length 7 to 30 that can contain
476 carbamidomethylated cysteine residues (fixed), oxidized methionine residues,
477 phosphorylated serine, threonine or tyrosine residues, Tandem Mass Tag (TMT)-labels,
478 TMTpro-labels and SILAC lysine4, arginine6, lysine8, and arginine10 stable isotope amino
479 acids that were generated by either collision-induced dissociation (resonance-type CID) or
480 higher-energy collisional dissociation (beam-type HCD). The model does not show any bias
481 in fragmentation prediction accuracy for tryptic vs. non-tryptic analytes (data not shown). The
482 INFERYS 3.0 retention time model can predict retention times for the same peptide classes
483 as above. The prediction capabilities of the deep learning models effectively dictate the
484 compatibilities of the CHIMERYS workflow.

485 **Brief description of the CHIMERYS algorithm**

486 **Setup**

487 The CHIMERYS workflow is a cloud-native API service, orchestrated by Kubernetes on
488 Amazon Web Services (AWS) or on-premise. The environment consists of two major
489 components: An INFERYS prediction server³⁷, which delivers predictions via gRPC requests

490 to a CHIMERYYS search algorithm instance, which matches these predictions to experimental
491 spectra.

492 **Description of the identification workflow**

493 The CHIMERYYS workflow follows the setup of classic search engines: After the *in-silico*
494 digest of the protein database and the generation of shuffled decoy sequences using a
495 similar logic as the mimic¹⁴ entrapment generator (see section below), a coarse first search is
496 performed to identify highly confident peptides for recalibration purposes. Notably, for a
497 group of I/L isomers, CHIMERYYS only scores one representative. A fast fragment ion index
498 implementation similar to MSFragger³⁸ is used to determine a ranked list of suitable
499 candidate peptides with isotope envelopes that (partially) overlap with the MS2 isolation
500 window (plus tolerances). Fragment ion intensities for highly-ranking candidate peptides are
501 predicted for each spectrum, merged against the experimental MS2 spectrum and
502 subsequently, a set of counting- (e.g. number of matching peaks between predicted and
503 experimental spectrum) and intensity-based scores (e.g. normalized spectral contrast angle)
504 are calculated. Candidate peptides that fall below certain cutoff criteria are removed. For
505 example, PSMs are required to have at least three matched fragment ions, one of which
506 must be the base peak (most abundant peak of the prediction) and another one of which
507 must be among the top three most intense peaks of the predicted spectrum. After the initial
508 search, a linear discriminant analysis identifies highly confident PSMs for the calibration of
509 optimal prediction parameters of the fragmentation model (e. g. normalized collision energy;
510 NCE), refinement learning of the retention time model and recalibration of fragment ion m/z
511 and match tolerances. Peptide classes with few confidently identified peptides are removed
512 entirely from the search space (e. g. peptides of length 7 carrying two missed cleavages and
513 two oxidized methionines). In the main search, the above-described scoring functions are
514 executed using the optimized settings and prediction parameters, albeit now also filtering
515 candidate peptides based on their predicted retention time. CHIMERYYS uses retention time
516 tolerance windows that would allow the identification of 99% of the peptides confidently
517 identified in the initial search. The scoring is repeated to arrive at a set of high-scoring
518 candidate peptides as input for the deconvolution function, where the candidate peptides
519 simultaneously compete for experimental fragment ion intensity in one concerted step.
520 CHIMERYYS uses non-negative L1-regularized regression via the LASSO, which models the
521 experimental spectrum as a function of the matrix of candidate peptides. The algorithm aims
522 to explain the spectrum with the fewest number of possible candidate peptides. This is
523 achieved by letting the different predicted spectra effectively 'compete' to explain the
524 experimentally observed fragment ion intensities. The algorithm reports a coefficient for each
525 candidate peptide representing its contribution to the experimental spectrum. A coefficient
526 bigger than 0 indicates that this candidate peptide was used to explain the experimental

527 spectrum. Based on the resulting coefficient, a subsequent round of intensity-based scoring
528 is executed. Here, the coefficients of the candidate peptides can be used to predict the
529 proportional intensity of all but one candidate peptide, add them together and subsequently
530 subtract this sum from the actual experimental MS2 spectrum to calculate what we call a
531 'shadow spectrum', i.e. the experimental spectrum with the contributions of all interfering
532 peptides removed. Next, the above-mentioned figures of merit are calculated based on these
533 shadow spectra without the interference of other peaks in the spectrum. Importantly, this also
534 works for fragment ion which are shared between candidate peptides. Candidate peptides
535 that fail to meet certain quality criteria (e. g. minimum number of most abundant peaks
536 shared between predicted and experimental spectrum) are filtered out. A list of all remaining
537 target and decoy PSMs per spectrum that received a coefficient > 0 and met all quality
538 criteria including all calculated scores is generated as input for the PSM-level error estimation
539 in mokapot¹³.

540 **FDR estimation using mokapot**

541 For error control, the initial implementation of CHIMERYYS utilized Percolator¹⁴ 3.0.5 to
542 aggregate all calculated scores for all target and decoy PSMs generated in a dataset. As
543 Percolator runtime scales poorly with large input files, we exchanged it with the Python-
544 based reimplementaion termed mokapot¹³. To ensure scalability to large input lists while
545 controlling the compute resources, we created a public pull request to the mokapot Github
546 repository (<https://github.com/wfondrie/mokapot/pull/100>) after we rewrote large parts of
547 mokapot's logic to allow streaming of data from disk, introduced RAM limits and implemented
548 more performant data structures. Mokapot is executed using the following parameters:
549 Training FDR of 1%, a training subset of 400k, and 10 iterations for training. We specifically
550 prevent mokapot from only retaining the top-scoring PSM per spectrum. Afterwards, the
551 resulting PSM-level q-values, support vector machine (SVM) scores and posterior error
552 probabilities (PEPs) are attached to the corresponding PSMs. Peptides containing
553 leucine/isoleucine (I/L) isomers in the search space are added back to the results with
554 identical scores and are flagged as *ambiguous*.

555 **MS2 quantification workflow**

556 CHIMERYYS determines raw file-specific peak apex retention times as the CHIMERYYS
557 coefficient-weighted mean of RT deltas relative to the gradient length based on PSMs
558 meeting 1% run-specific PSM-level FDR. If an external inclusion file was used, PSMs
559 meeting 1% run-specific PSM-level FDR including their relative retention times and
560 CHIMERYYS coefficients from the list are also considered. If no PSMs meet 1% run-specific
561 PSM-level FDR for a given precursor in a given raw file, the apex for said precursor in this
562 raw file is calculated using the same logic as above, but based on PSMs meeting 1% run-

563 specific PSM-level FDR in other raw files and the inclusion file. CHIMERYYS in its current
564 implementation then estimates maximum integration borders per raw file as the 99% quantile
565 of peak widths at base (not full width at half maximum; FWHM) from precursors with at least
566 three PSMs surviving a run-specific PSM-level FDR threshold of 1%. These maximum
567 integration borders are then applied to each precursor in this raw file, leading to relatively
568 wide integration borders, particularly for low-abundant precursors. Afterwards, quantification
569 of PRM and DIA data is performed by either trapezoidal integration of the CHIMERYYS
570 coefficients from each precursor in a set of consecutive MS2 spectra sharing the same
571 isolation window within the integration borders, or by using the highest CHIMERYYS
572 coefficient within the integration borders as the elution peak apex intensity. One missing
573 CHIMERYYS coefficient in a series of consecutive MS2 scans with the same isolation window
574 is allowed (gap scan) and a contribution of 0 is inserted to any further scan with missing data
575 points, which act as boundary for peak area integration. Notably, at this point, CHIMERYYS
576 coefficients are taken from PSMs irrespective of their run-specific PSM-level FDR. However,
577 CHIMERYYS coefficients will only be used from peptide precursors that met CHIMERYYS'
578 quality criteria (e. g. a minimum of three peaks matched between the predicted and the
579 experimental spectrum) and are located in the vicinity of the determined peak apex. As such,
580 at least one confidently identified PSM across all raw files is required to generate quantitative
581 values based on PSMs around the determined peak apex in each raw file. As such,
582 CHIMERYYS will quantify precursors that fail to meet run-specific precursor-level FDR
583 thresholds. Users are free to filter their list of precursors at 1% global precursor-level FDR
584 (precursor was confidently identified in at least one raw file) or additionally also at 1% run-
585 specific precursor-level FDR. The latter will reduce data completeness and is more
586 conservative. However, we have shown that often, these quantifications are precise and
587 accurate, so we recommend to work with precursors filtered to 1% global precursor-level
588 FDR during exploratory data analysis and turn to run-specific precursor-level FDR for the
589 validation of interesting hits.

590 **Post-processing of CHIMERYYS' PSM-level outputs**

591 CHIMERYYS 2.7.9 is integrated into Thermo Fisher Scientific Proteome Discoverer software
592 3.1 (PD)²⁷. Hence, PD starts CHIMERYYS searches on AWS by uploading an internal format
593 containing only MS2 spectra and some auxiliary information, a fasta file and the search
594 parameters to the CHIMERYYS service, which then processes the data and generates a result
595 file. This result file is then downloaded and post-processed by PD²⁷. In this study, we used
596 the default CHIMERYYS processing and consensus workflows with minor modifications.
597 Briefly, all DDA data processing was using the PSM-grouper node to generate peptide
598 groups, which were then validated using quality³⁹. For DIA data, we used a special PCM-
599 grouper node, which enables the calculation of run-specific and global precursor-level FDR.

600 MS1-based quantification was performed using the Minora Feature Detector with default
601 settings. MS2-based quantification was performed using the MS2 Fragment Ions Quantifier
602 node with default settings.

603 **Data generation**

604 **External data**

605 The following external data were downloaded from PRIDE and processed with the respective
606 search engines. In brief, ASTRAL data extracted from Gutzman et al., 2024 (PXD046453)²⁹,
607 Arabidopsis and Halobacterium data from Müller et al., 2020 (PXD014877)⁴⁰, body fluid data
608 from Bian et al., 2020 (PXD015087)⁴¹, secretome data from Tüshaus et al., 2020
609 (PXD018171)⁴² and triple species mix as well as HeLa data from the LFQBench-type dataset
610 by Van Puyvelde et al., 2022 (PXD028735)²². Notably, peptides carrying oxidized methionine
611 residues were excluded from all analyses of the LFQBench-type dataset, since raw files
612 showed differential oxidation (data not shown). An itemized mapping of external data
613 processed as part of this study to their source and respective search settings will be made
614 available upon publication.

615 **Internal data**

616 The following datasets were generated in house: FFPE, gradient comparison, wwDDA,
617 instrument generations and PRM data.

618 **Cell culture and sample preparation**

619 Human HeLa and pancreatic mouse cells (ATCC, CCL-2) were cultured under standard
620 conditions at 37°C with 5% CO₂ in DMEM medium supplemented with 10% fetal bovine
621 serum (FBS) and 100 U/mL penicillin (Invitrogen). At around 80% confluence cells were
622 washed three times with PBS buffer before UREA lysis (8M Urea, 80 mM Tris pH 7.6, 1x
623 protease inhibitor) was performed for 5 minutes on ice. Cell lysate was clarified by
624 centrifugation (20,000 x g for 10 min).

625 In solution protein digest was conducted as following. First, proteins were reduced with 10
626 mM DTT at 37°C for 1 h, followed by alkylation with 2-chloroacetamide (CAA) at a final
627 concentration of 55 mM for 45 min at room temperature in the dark while shaking on a
628 thermo shaker. After the addition of five volumes of 50 mM Tris (pH 8), trypsin digest was
629 performed overnight by adding trypsin twice (1:100) after a primary incubation time of 4h.
630 Desalting was performed using Sep-Pak columns according to the user manual. Human
631 brain FFPE samples were digested using a SDS lysis protocol followed by digestion with the
632 SP3 approach as described in detail in Tüshaus et al. 2023⁴².

633 **LC-MS/MS**

634 FFPE, gradient comparison and wwDDA data were acquired on a micro-flow LC coupled via
635 a HESI source to an Q Exactive HF-X hybrid quadrupole-Orbitrap mass spectrometer
636 (Thermo Scientific). Optimization of the micro-flow LC setup as well as technical details were
637 previously published in Bian et al., 2020⁴¹. In brief, peptide separation was performed on an
638 Acclaim™ PepMap™ 100 C18-HPLC-column (15 cm lengths, 1 mm inner diameter, 2 μm
639 particle size, #164711, Thermo Fisher Scientific) at 55°C. Linear gradients with buffer A
640 (0.1% FA, 3% DMSO in dH₂O) and buffer B (0.1% FA, 3% DMSO in ACN) from 3 to 28% B
641 were run at 50 μL/min. Sample loading, column wash and equilibration was performed at 100
642 μL/min. Source settings were applied as following: 320°C capillary temperature, 3.5 kV spray
643 voltage, 300°C auxiliary gas. MS data were acquired at a normalized collision energy of 28
644 %, in Top20 mode, at an m/z range of 360–1300, AGC target of 3E6 and 1E5, maximal
645 injection time of 50 ms and 22 ms, resolution of 60 k and 15 k on MS1 and MS2 level,
646 respectively. The MS2 isolation window width was 1.4 Th in standard DDA runs and
647 increased up to 20.4 Th for wide window acquisition DDA as indicated in the figure legends.

648 Ion trap data were acquired with an Orbitrap Eclipse™ Tribrid™ mass spectrometer (Thermo
649 Scientific) that was coupled to a Dionex UltiMate 3000 RSLCnano System (Thermo
650 Scientific). Samples were transferred onto a trap column (75 μm x 2 cm, 5 μm C18 resin
651 Reprosil PUR AQ - Dr. Maisch). After washing with the trap washing solvent (5 μl/min, 10
652 min), samples were separated on an analytical column (75 μm x 48 cm, 3 μm C18 resin
653 Reprosil PUR AQ - Dr. Maisch). A 70 min method, including a 50 min gradient, was
654 performed with a flow rate of 300 nl/min (4 % B up to 32 % B within 50 min). Solvent A: 0.1
655 %v FA, 5%v DMSO in dH₂O, Solvent B 0.1 %v FA, 5 %v DMSO in ACN. MS1 scans were
656 acquired with an Orbitrap resolution of 60 k, within a scan range of 360-1300 m/z, a
657 maximum injection time of 50 ms, a normalized AGC target of 100% and RF lens of 40 %,
658 including charge stage 2 to 6, with an exclusion time of 25 s. MS2 scans were performed with
659 the ion trap with a normalized AGC target of 200%, a maximum injection time of 25 ms and
660 either with an HCD collision energy of 31% (wwDDA) or with an CID collision energy of 35 %
661 (CID). Quadrupole isolation window was varied between 0.4, up to 5.0 m/z as indicated in the
662 figure.

663 Data for the instrument comparison were assembled from 1-hour HeLa QC runs, acquired
664 over several years at the Chair of Proteomics and Bioanalytics at the TUM. They were run on
665 various liquid chromatography systems, employed diverse instrument-specific settings,
666 slightly different gradients and used different batches of HeLa digest, prepared in-house.

667 **Targeted assay generation**

668 A simple PRM assay was devised by randomly selecting 18 proteins and 2-3 peptide
669 precursors each across the whole measured intensity range from a 1-hour HeLa run
670 analyzed on a Orbitrap Fusion Lumos mass spectrometer (Thermo Scientific). A total of 51
671 precursors were put into an inclusion list in addition to 14 precursors corresponding to a
672 retention time standard. A 1-hour HeLa sample was analyzed in PRM mode: MS2 spectra
673 were acquired using 0.4 Th isolation window, a maximum injection time of 100 ms, HCD
674 collision at a normalized collision energy of 28 % and read out in the Orbitrap at 15k
675 resolution.

676 **Data processing and evaluation**

677 **Data analysis using various search algorithms**

678 A variety of search engines have been applied in the scope of this study. Unless otherwise
679 noted, default settings were used including trypsin digestion (with proline rule, so cutting after
680 K/R, but not before P), carbamidomethylation of cysteine as fixed modification as well as
681 methionine oxidation as variable modification. Peptide lengths were limited to 7-30 for
682 comparisons between CHIMERYS, DIA-NN and Spectronaut. All searches were performed
683 against the canonical FASTA files of the corresponding species or species mix. A
684 contaminants FASTA was utilized in all searches to control for contaminations⁴³. Entrapment
685 searches were performed as described in detail its own section. All FASTA files used in this
686 study will be made publicly available via PRIDE upon publication.

687 **CHIMERYS.** Searches were performed using CHIMERYS v2.7.9 from PD v3.1.0.622 using
688 default settings. For quantification, non-normalized intensities from Minora Feature Detector
689 were used. PSM, precursor, peptide group and protein group IDs were extracted from
690 .pdResult files after the removal of decoys and PSM-, precursor-, peptide group- or protein-
691 level q-value filtering. For ion trap data, an additional TopN Peaks filter node was used
692 before the CHIMERYS node with TopN set to 15 and mass window to 100 Da.

693 **Sequest HT.** Searches were performed from PD v3.1.0.622 using default settings. For
694 quantification, non-normalized intensities from Minora Feature Detector were used. PSM and
695 peptide group IDs were extracted from .pdResult files after the removal of decoys and PSM-
696 or peptide group-level q-value filtering.

697 **MS Amanda.** Searches were performed using MS Amanda v2.0 from PD v3.0.0.201 using
698 default settings. For quantification, non-normalized intensities from Minora Feature Detector
699 were used. PSM and peptide group IDs were extracted from .pdResult files after the removal
700 of decoys and PSM- or peptide group-level q-value filtering.

701 **Comet.** Searches were performed from PD v3.1.0.622 using default settings. For
702 quantification, non-normalized intensities from Minora Feature Detector were used. PSM and
703 peptide group IDs were extracted from .pdResult files after the removal of decoys and PSM-
704 or peptide group-level q-value filtering.

705 **MSFragger.** Searches were performed using MSFragger v20.0 using the "Default" workflow.
706 Additionally, MSFragger v21.1 was utilized with the "WWA" workflow and "DDA+" data type,
707 which considers candidate peptides in the full MS1 isolation window and reports up to the top
708 five PSMs instead of only one. Precursor-level IDs were extracted from ion.tsv files, which
709 were already filtered for decoys and precursor FDR. Peptide group-level IDs were rolled up
710 from precursor-level data, controlled at 1% precursor-level FDR.

711 **MetaMorpheus.** Searches were performed using MetaMorpheus v0.0.320 with default
712 settings. PSM-level IDs were extracted from AllPSMs.psmtsv files after the removal of
713 decoys and PSM-level q-value filtering. MetaMorpheus does not feature peptide group-level
714 IDs, so these were rolled up from the PSM level, controlled at 1% PSM-level FDR.

715 **MS-GF+.** Searches were performed using MS-GF+ v2022.04.18 with default settings. The
716 output .mzid file was converted to .tsv, from which PSM-level IDs were extracted after the
717 removal of decoys and PSM-level q-value filtering. MS-GF+ does not feature peptide group-
718 level IDs, so these were rolled up from the PSM level, controlled at 1% PSM-level FDR.

719 **MaxQuant.** Searches were performed with default settings in MaxQuant v2.4.2.0. Peptide
720 group IDs were extracted from modificationSpecificPeptides.txt files after filtering for decoys,
721 which were already peptide group-level FDR filtered.

722 **DIA-NN.** Library-free searches were performed using DIA-NN 1.8.1²³. Mixed species samples
723 were analyzed with a classic, concatenated, and peptide eFDR approach. Entrapment fastas
724 were based on a database for the three species and a database containing contaminants.
725 Detailed information on entrapment database generation can be found under 'Entrapment
726 database construction via mimic'. For all searches, spectral libraries were generated from
727 fasta files using DIA-NN's prediction capabilities. Default settings were used with the
728 following adjustments: Missed cleavages = 0, Maximum number of variable modifications =
729 2, Quantification strategy = Peak height. 'Ox(M)' was added as a variable modification. Note
730 that 'MBR' is checked per default for library-free searches. The output results from the
731 report.tsv were used. For run-specific precursor filtering, results were filtered for Q.Value ≤
732 0.01. Precursor.Quantity was used as precursor quantification values. Fragments for
733 identification were counted considering fragments with Fragment.Quant.Corrected > 0.
734 Fragments for quantification were counted by identifying fragments, for which the sum of

735 Fragment.Quant.Corrected corresponded to Precursor.Quantity. Samples with different
736 species composition were analyzed together.

737 **Spectronaut.** Raw files were analyzed using Spectronaut v18 (Biognosys) with a library-free
738 approach (directDIA+). Mixed species samples were analyzed with a classic, concatenated,
739 and peptide eFDR approach. Entrapment fastas were based on a database for the three
740 species and a database containing contaminants. Detailed information on entrapment
741 database generation can be found under 'Entrapment database construction via mimic'.
742 Default settings were used with the following adjustments: Max Peptide Length = 30, Missed
743 Cleavages = 0, Max Variable Modifications = 2, Quantity Type = Height. Carbamidomethyl
744 (C) was set as a fixed modification, and Oxidation (M) was set as a variable modification.
745 Cross-run normalization was disabled. Samples with different species composition were
746 analyzed together. For run-specific precursor filtering, results were filtered for EG.Qvalue \leq
747 0.01. `EG.TotalQuantity (Settings)` was used as precursor quantification values. Fragments
748 for identification were counted considering fragments with F.PeakArea > 1. Fragments for
749 quantification were counted by identifying fragments with F.PeakArea > 1 and
750 F.ExcludedFromQuantification = FALSE.

751 **Post-hoc filtering of Spectronaut output files**

752 We show that Spectronaut claims to confidently identify precursors with no apparent signal in
753 the corresponding MS2 spectra. In order to remove these precursors, an unfiltered, fragment-
754 level export needs to be created from within Spectronaut. Based on this report, fragments
755 with F.PeakArea > 1 and F.ExcludedFromQuantification = FALSE can be identified. In
756 Spectronaut, the desired number of fragments for quantification can be set (under DIA
757 Analysis > Quantification > Interference Correction > MS2 Min). Per default, 'MS2 Min' is
758 three. Hence, in this study, precursors were filtered to only those with at least 'MS2 Min'
759 fragments that have F.PeakArea > 1 and F.ExcludedFromQuantification = FALSE.

760 **Calculation of shared fragment ions**

761 Peptide fragment intensity predictions of PSMs identified by CHIMERYS v2.7.9 were
762 generated using INFERYS 3.0.0. The R package rawrr⁴⁴ was used to extract centroided m/z
763 and intensity values from raw files. Different PSMs within the same MS2 scans were tested
764 for shared fragment ions using two methods: directly matching predicted fragment ions (raw
765 file independent) and checking if predicted fragment ions matched to the same raw file peaks
766 with a 20 ppm tolerance. The fraction of shared fragment ions among all fragment ions was
767 determined and analyzed based on amino acid position and 200 m/z bins.

768 **Entrapment database construction via mimic**

769 We generated same-organism entrapment proteins in nine-fold excess relative to the number
770 of target proteins by shuffling sequences at the peptide level using mimic¹⁴ with the command
771 line flags --prepend, --empiric, --replacel and --mult-factor 9. The goal is to generate nine
772 different entrapment peptides per target peptide, which are isobaric, shuffled versions of it.
773 Entrapment peptides shall have the same peptide N- and C-termini as their target
774 counterpart and shall be as close in amino acid composition to it as possible. Briefly, mimic
775 achieves this by first digesting the target protein database into fully tryptic peptides without
776 missed cleavages (no proline rule, i. e. cleaving proteins after each K/R). Afterwards, the
777 sequence of each unique target peptide is permuted randomly while keeping the termini
778 fixed. In case permutation of a peptide generates a target or previously-generated
779 entrapment peptide, shuffling is repeated 1,000 times. For this comparison, isoleucine and
780 leucine are considered to be identical amino acids. If no entrapment peptide could be
781 generated after 1,000 rounds of shuffling, amino acids are mutated. Briefly, a random amino
782 acid is mutated to a different one while respecting the amino acid frequencies of the target
783 protein database. Isoleucine and leucine are never mutated when using the --replacel flag.
784 Notably, it is ensured that target peptides with the same peptide sequence generate the
785 same entrapment peptides. Subsequently, entrapment peptides are assembled to
786 entrapment proteins, which are then either appended to the target database in the case of
787 the classic fasta concatenation approach (classic eFDR) or concatenated to the
788 corresponding target protein sequences – separated by lysine residues – in the sequence
789 concatenation approach (concatenated eFDR). The final digested fasta approach (peptide
790 eFDR) used a peptide-level fasta, which was generated by digesting the classic eFDR fasta
791 file with Protein Digestion Simulator v2.4.7993.32903 before passing it to the search engine
792 (fully tryptic digest with no missed cleavages and no proline rule). Notably, protein-C-terminal
793 peptides cannot be identified with the concatenated eFDR approach and neither can protein
794 N-terminal peptides with methionine excision from entrapment proteins. Mimic is available at
795 <https://github.com/percolator/mimic>. A simplified web-version of it is available at
796 <https://mimicerys.msaid.io/>.

797 **FDR definitions**

798 We adhere to previously-established definitions of FDR in the run-specific and global context
799 for any identification level (e. g. precursors, modified peptides, peptides, peptide groups,
800 protein groups)²⁵. For example, peptide-level FDR in the run-specific context can answer the
801 question “Which peptides were detected at 1% FDR in this specific LC-MS/MS run?”.
802 Peptide-level FDR in the global context can instead answer the question “Which peptides

803 were detected at 1% FDR in at least one LC-MS/MS run of a given experiment containing
804 multiple LC-MS/MS runs?”.

805 **Calculation of entrapment FDR**

806 We empirically validated the self-reported FDR by CHIMERY5, DIA-NN and Spectronaut
807 using entrapment experiments, which are sometimes also called double-decoy experiments.
808 The general idea is that most search engines are black boxes that report lists of
809 identifications at self-reported FDR. Entrapment experiments try to empirically validate this
810 self-reported FDR.

811 Briefly, most if not all search engines try to control the FDR of their identifications at some
812 percentage (usually 1 %). In practical terms, this means they try to limit for example the
813 number of random false target PSMs in the final list of PSMs, i. e. the FDR at the PSM level.
814 However, because the identity of random false target identifications is not known *a priori*,
815 search engines model the behavior of random false targets with the help of decoys. Decoys
816 are peptides that resemble target peptides, but are assumed to be absent from the sample
817 under investigation. Many if not all search engines generate their own decoy peptides using
818 one of several approaches (e. g. reversing protein sequences⁸). Experimental MS2 spectra
819 are then compared to theoretical, predicted or library spectra from target and decoy peptides,
820 which results in a score for each PSM that measures how closely they match. In classic
821 target decoy competition, only the top-scoring PSM of each spectrum is retained. The
822 remaining PSMs are then sorted in descending order by their score (assuming that a high
823 score is associated with a good match) and q-values are calculated as the cumulative sum of
824 decoys divided by the cumulative sum of targets for each PSM. Usually, these q-values are
825 then monotonized by sorting the PSMs by their score in ascending order and calculating the
826 cumulative minimum of the q-values for each PSM to arrive at the final q-values. Removing
827 PSMs with a q-value > 0.01 would filter the list of PSMs to 1% FDR. Notably, as mentioned
828 above, this approach assumes that decoys resemble random false targets and consequently
829 that the score distribution of decoy peptides is identical in shape and magnitude to the score
830 distribution of random false target peptides. If this is not the case, FDR could be over- or
831 underestimated.

832 Entrapment or double-decoy experiments measure empirical FDR by following very similar
833 principles as the target/decoy approach described above. Before passing the target protein
834 database to the search engine, which then generates its own decoys (one for every peptide
835 after *in silico* digestion of the target protein database), so-called entrapment proteins are
836 added to it. Entrapment proteins consist of entrapment peptides, which are very similar to
837 decoy peptides in that they should resemble target peptides, but are assumed to be absent
838 from the sample under investigation. However, the search engine does not know which

839 peptide is a target peptide and which peptide is an entrapment peptide. As such, it generates
840 decoy peptides for each target and each entrapment peptide. We will call them 'normal
841 decoys' and 'entrapment decoys', respectively. Experimental MS2 spectra are then
842 compared to theoretical, predicted or library spectra from target, normal decoy, entrapment
843 and entrapment decoy peptides. Now the search engine calculates its FDR as described
844 above after sorting PSMs based on their score in descending order by dividing the
845 cumulative sum of decoy identifications (normal decoys + entrapment decoys) by the
846 cumulative sum of target identifications (targets + entrapments), followed by q-value
847 monotonicity. In addition, since the identity of entrapments is known to the researcher, an
848 entrapment FDR can be calculated. This is done by sorting PSMs in descending order by
849 their score and calculating entrapment q-values as the cumulative sum of normal decoys plus
850 entrapments, divided by the cumulative sum of targets plus entrapments for each PSM,
851 followed by entrapment q-value monotonicity. For entrapment analyses with DIA-NN and
852 Spectronaut, the formula for calculating entrapment FDR was slightly modified, because their
853 final results violate the assumption that the score distribution of decoy peptides is identical in
854 shape and magnitude to the score distribution of random false target peptides. This is
855 because they perform two searches of the data; one to generate a filtered spectral library and
856 a second one that uses this spectral library to analyze the same data again in a peptide-
857 centric fashion. The filtering of the spectral library is done based on an initial FDR calculation,
858 which will remove a substantial fraction of random false target peptides, causing the
859 distribution of decoy peptides in the second search to not be identical anymore in shape and
860 magnitude to the score distribution of random false target peptides. Hence, the cumulative
861 sum of normal decoys was removed from the entrapment FDR calculation for DIA-NN and
862 Spectronaut. As such, entrapment FDR was calculated by sorting precursors in descending
863 order by their score and calculating entrapment q-values as the cumulative sum of
864 entrapments, divided by the cumulative sum of targets plus entrapments for each precursor,
865 followed by entrapment q-value monotonicity. Removing PSMs with an entrapment q-
866 value > 0.01 should also filter the list of PSMs or precursors to 1% FDR if the search engine
867 does not have a bias in its scoring function. In other words, if there is no bias, then q-values
868 and entrapment q-values should follow the diagonal if they are plotted against one another
869 in a scatter plot.

870 Given the entrapment FDR calculations above, it is clear why we use entrapments in at least
871 nine-fold excess relative to targets. The reason is because the search engine generates
872 decoys for each target and each entrapment. In other words, using the same number of
873 entrapments as targets would result in a search space consisting of one part targets, two
874 parts decoys (one part normal decoys and one part entrapment decoys) and one part
875 entrapments. Given the nature of random false identifications, they are equally likely to map

876 to false targets, decoys or entrapments. However, we cannot identify them as false matches
877 if they map to false targets. Consequently, the likelihood of a random false identification that
878 we can identify to map to decoys is twice as high (66%) as the likelihood of it mapping to an
879 entrapment (33%) if the number of targets and entrapments was the same. As we increase
880 the number of entrapments that are added to the target protein database, the likelihood of a
881 random false identification to map to decoys or entrapments becomes more and more
882 similar. At a nine-fold excess of entrapments relative to targets, the likelihood of a random
883 false identification that we can identify to map to a decoy is merely ~6% higher (53%) than
884 the likelihood of it mapping to an entrapment (47%).

885 Entrapments can be added to the target protein database in various different ways (see
886 section 'Entrapment database construction via mimic'). In the case of the concatenated
887 eFDR approach, we had to modify the calculation of eFDR slightly for CHIMERYYS, because
888 it is no longer possible to distinguish between normal decoys and entrapment decoys. This is
889 because they are generated internally by CHIMERYYS and are both annotated with the target
890 protein identifier, which can be used in the other entrapment approaches to identify normal
891 decoys and entrapment decoys. Hence, for the concatenated eFDR approach, eFDR
892 calculation is done by sorting PSMs in descending order by their score and calculating
893 entrapment q-values as the cumulative sum of decoys plus entrapments, divided by the
894 cumulative sum of targets plus entrapments for each PSM, followed by entrapment q-value
895 monotonicization. No change was made to the calculation of eFDR based on the concatenated
896 eFDR approach for DIA-NN and Spectronaut.

897 Throughout the manuscript, we performed several entrapment analyses for CHIMERYYS, DIA-
898 NN and Spectronaut. In Figure 1B and Supplementary Figure 8, we directly compared the
899 self-reported run-specific FDR at a given identification level (e. g. peptides or precursors) to
900 the corresponding entrapment FDR. In Figures 3 and 4, as well as in Supplementary Figure
901 9, we filtered identifications either at the self-reported run-specific FDR or additionally also at
902 the corresponding entrapment FDR. We then compared the characteristics of the
903 identifications that only survive the self-reported FDR or also survive the entrapment FDR (e.
904 g. precision and accuracy of quantification). For CHIMERYYS, we used the q-values and SVM
905 scores reported by qvality in PD, which either validated peptide groups in the global context
906 (Figure 1B) or precursors in the run-specific context (all other visualizations). For
907 Spectronaut, we used the EG.Qvalue and the EG.Cscore columns of unfiltered Spectronaut
908 exports, which validated precursors in the run-specific context. For DIA-NN, we used the
909 Q.Value and CScore columns of the main report.tsv, which validated precursors in the run-
910 specific context. Unfortunately, DIA-NN and Spectronaut do not report q-values for decoy
911 identifications and the former always filters its report.tsv to 1% run-specific precursor-level
912 FDR. They do however report scores for decoy identifications. In order to convert these

913 scores into q-values for the plots in Supplementary Figure 8, we interpolated the q-values
914 given the relationship of score and q-value for target precursors in each raw file separately.

915 **PRM Analysis in CHIMERYs and Skyline**

916 PRM data were searched using CHIMERYs against a Human Swissport database with
917 default settings. The resulting MS2-based quantities and coefficient traces were extracted
918 from the .pdresult files using custom R scripts. The .raw file was also analyzed in Skyline⁹
919 v22.2 using a ProSight-predicted spectral library¹⁵ as integrated in Skyline. Peak boundaries
920 were manually refined for all precursors and the quantitative values for the area beneath the
921 five most abundant fragment ions was aggregated into a quantitative measure. Correlations
922 were calculated and visualized using custom R scripts.

923 **Comparison of DDA and DIA data using CHIMERYs**

924 Unless otherwise noted, DDA and DIA files were searched with CHIMERYs v2.7.9 from PD
925 v3.1.0.622, using Minora for MS1-level and CHIMERYs for MS2-level quantification,
926 respectively. PSMs were filtered at 1% run-specific FDR, peptide groups and protein groups
927 at 1% dataset global peptide group and protein FDR, respectively. Peptide groups and
928 protein groups with FoundInSamples = 0 were discarded. Protein groups were further filtered
929 for PsmCount > 0 and IsMasterProtein = 0. For quantification of conditions / ratios, at least 2
930 non-normalized intensity values > 0 were required per condition / in both conditions,
931 respectively (or 3 if stated “all replicates”). For calculating CVs, at least 3 non-normalized
932 intensity values > 0 were required per condition. 2D density estimates were calculated using
933 kde2d from the R package MASS. Entrapment FDR was calculated as described above
934 using the peptide eFDR approach.

935

936 Data Availability

937 External raw data

938 An itemized mapping of external data processed as part of this study to their source and
939 respective search settings will be made available upon publication.

940 Generated data

941 The generated mass spectrometric raw and search data of internal datasets from this study
942 will be made available via PRIDE⁴⁵ upon publication.

943 Open-source software

944 The mokapot version used in this study is available as a public pull request on GitHub
945 (<https://github.com/wfondrie/mokapot/pull/100>). The modifications to the mimic entrapment
946 database generator are available on GitHub (<https://github.com/percolator/mimic/>). A web-
947 version of the mimic tool can be found at <https://mimicerys.msaid.io/>.

948 Data analysis and plotting scripts

949 The custom R scripts underlying the data analysis are available upon request.

950 Acknowledgements

951 The authors wish to thank numerous scientific colleagues for their input, discussions, and
952 support. The authors want to expressively thank the Proteome Discoverer (PD) software
953 development team at Thermo Fisher Scientific for their collaboration, support, and
954 contributions on the successful integration of CHIMERY5 into PD and the scientific discourse
955 on the results. The authors wish to thank Elmar Zander for consulting on mathematical
956 topics. The authors wish to thank Matthew The for consulting on entrapment experiments
957 and FDR control. The authors also wish to thank their colleagues Dulguun Bold, Jeremiah
958 Santoso, Agnes Guevende and Mohammed Al Kiddeh for various contributions to the
959 software.

960 Author contributions

961 M.F. and M.W. conceived the study. M.F., M.W., and D.P.Z. developed and evaluated the
962 initial prototype. F.S., P.S., T.S., M.G., I.B., S.B.F. and S.G. developed, implemented, and
963 optimized the algorithms. S.G., V.S., S.B.F., L.M. and M.G. developed the deep learning
964 models. T.S., M.G. and F.S. orchestrated the implementation of software modules and the
965 deployment of the software. M.F., D.P.Z., F.S., M.G., M.T.B. and A.H. evaluated the
966 algorithm. M.F., M.T.B., J.T., A.H., and D.P.Z. processed the result data. M.F., M.T.B., J.T.,
967 A.H. and D.P.Z. performed the data analysis. L.E. helped in the preparation of the Figures.
968 M.F., D.P.Z., M.B., A.H., F.S., P.S., S.G., T.S., J.T., B.K. and M.W. provided critical

969 feedback, discussed the results, and consulted in revisions. M.F., D.P.Z., J.T., B.K. and M.W.
970 wrote the manuscript.

971 **Competing interest statement**

972 M.F., D.P.Z., S.G. and T.S. are co-founders, shareholders, and employees of MSAID GmbH,
973 a company that develops software for proteomics, including the algorithm presented in this
974 manuscript. M.W. and B.K. are co-founders and shareholders of MSAID GmbH and
975 OmicsScouts GmbH, which operates in the field of proteomics, but they have no operational
976 role in either company. M.T.B., A.H., F.S., M.G., P.S., S.B.F., V.S., L.E., I.B., L.M., and M.S.
977 are employees of MSAID GmbH.

978

979 Figure Legends

980 Figure 1 – Deconvolution of chimeric DDA spectra:

981 **(A)** CHIMERYYS treats chimeric spectra as linear combinations of pure spectra, and its
982 spectrum-centric deconvolution uses non-negative regularized regression to estimate
983 interference-corrected total ion currents for all candidate peptides. **(B)** Peptide-level
984 entrapment analysis with the classic eFDR approach (see Online Methods) of DDA data
985 acquired using different isolation windows widths and processed with CHIMERYYS. **(C)**
986 Example of a deconvoluted chimeric spectrum with six PSMs from a 2-hour HeLa DDA
987 single-shot measurement, acquired on an Orbitrap QE HF-X with 1.3 Th isolation windows
988 from the LFQbench-type dataset²². The inset visualizes how the experimental intensity of
989 shared fragment ions with low m/z values is distributed to multiple PSMs. **(D)** Peptide-level
990 density plots for triplicate 2-hour DDA single-shot measurement from two different conditions,
991 acquired on an Orbitrap QE HF-X with 1.3 Th isolation windows from the LFQbench-type
992 dataset, analyzed with CHIMERYYS (top) or Sequest HT (bottom). FDR was controlled at 1%
993 at the global peptide group-level. **(E)** Comparison of peptide identifications from multiple
994 search engines on the same data as in **(C)**. FDR was natively controlled at different levels,
995 depending on the search engine.

996 Figure 2 – Optimizing data acquisition with deconvolution in mind

997 **(A)** Peptide identifications at 1% global peptide group-level FDR based on Sequest HT
998 (orange) and CHIMERYYS (blue) from 1h HeLa single-shot measurements, acquired using
999 various Orbitrap generations. **(B)** PSM, peptide and protein group identifications based on
1000 Sequest HT (orange), CHIMERYYS (dark blue) and CHIMERYYS after removal of low-
1001 abundance peaks (light blue) from a 1-hour HeLa single-shot measurement, acquired using
1002 CID fragmentation with ion trap read out. FDR was controlled at 1% at the run-specific PSM-,
1003 global peptide group- and global protein level, respectively. **(C)** PSM, peptide and protein
1004 group identifications based on Sequest HT (orange) and CHIMERYYS (blue) from HeLa
1005 single-shot measurements using various gradient lengths. FDR was controlled at 1% at the
1006 run-specific PSM-, global peptide group- and global protein level, respectively. **(D)**
1007 Distribution of the number of PSMs per MS2 spectrum from 1-hour HeLa single-shot
1008 measurements, acquired using different isolation window widths. FDR was controlled at 1%
1009 at the run-specific PSM-level. **(E)** PSM, peptide and protein group identifications based on
1010 Sequest HT (orange) and CHIMERYYS (blue) from 1-hour HeLa single-shot measurements,
1011 acquired using different isolation window widths. FDR was controlled at 1% at the run-
1012 specific PSM-, global peptide group- and global protein level, respectively.

1013 **Figure 3 – Deconvolution of chimeric DIA spectra**

1014 **(A)** Precursors quantified by CHIMERYYS, DIA-NN and Spectronaut in at least one (orange)
1015 or two (gray) out of three replicate measurements of two different conditions in triplicate 2-
1016 hour DIA single-shot measurement from two different conditions, acquired on an Orbitrap QE
1017 HF-X with 8 Th isolation windows from the LFQbench-type dataset. Identifications are filtered
1018 at 1% run-specific precursor-level FDR or additionally also at 1% run-specific precursor-level
1019 eFDR (Online Methods). **(B)** Apex intensities for precursors surviving (gray) or not surviving
1020 (red) 1% run-specific precursor-level eFDR for the same data as in **(A)**, analyzed by
1021 CHIMERYYS, DIA-NN and Spectronaut. **(C)** Number of fragment ions used for the
1022 quantification of precursors by CHIMERYYS, DIA-NN and Spectronaut for the same data as in
1023 **(A)**. Identifications are filtered at 1% run-specific precursor-level FDR and colored by whether
1024 they also survive 1% run-specific precursor-level eFDR (Online Methods). **(D)** Example
1025 fragment ion XICs for two high-scoring precursors identified by Spectronaut (all six library
1026 fragments are shown). The top one is also identified by CHIMERYYS, while the bottom one is
1027 not. The vertical dashed gray lines mark Spectronaut's reported EG.StartRT and EG.EndRT
1028 for the corresponding precursors, with 102.36 min and 102.94 min for the top and 19.75 min
1029 and 20.05 min for the bottom XIC, respectively.

1030 **Figure 4 – Coefficient-based quantification**

1031 **(A)** CHIMERYYS quantifies precursors based on MS2 spectra by tracing their spectrum-
1032 centric coefficients – interpretable as interference-corrected TIC – over retention time,
1033 followed by apex intensity extraction or trapezoidal peak area approximation. **(B)** Correlation
1034 of CHIMERYYS' MS2-based quantification to Skyline's MS2-based quantification on a PRM
1035 dataset targeting 52 peptides from 18 human proteins across the entire dynamic range of a
1036 human proteome. **(C)** Violin plots depicting CV distributions of precursors identified by
1037 CHIMERYYS, DIA-NN or Spectronaut at 1% self-reported FDR (top) or additionally requiring at
1038 least 1 replicate per condition to also meet 1% entrapment FDR (bottom, peptide eFDR
1039 approach) at the precursor-level in the run-specific context, quantified based on 0, 1, 2, and 3
1040 or more fragment ions. Data is from triplicate 2-hour DIA single-shot measurement from two
1041 different conditions, acquired on an Orbitrap QE HF-X with 8 Th isolation windows from the
1042 LFQbench-type dataset. Vertical dashed lines depict the median across all search engines
1043 with the respective filtering applied (29.6% for filtering at self-reported FDR and 28.9% with
1044 additional eFDR filtering). **(D)** Precursor-level log₂-ratio density plots for the same data as in
1045 **(C)**, stratified by whether at least one replicate per condition survived 1% eFDR based on the
1046 peptide eFDR approach (top) or not (bottom). Replicate measurements were averaged
1047 before calculating ratios.

1048 **Figure 5 – Comparing DDA and DIA data**

1049 **(A)** Peptide groups identified by CHIMERYs in triplicate 2-hour DDA (left) or DIA (right)
1050 single-shot measurement from two different conditions from the LFQbench-type dataset,
1051 acquired on an Orbitrap QE HF-X with 1.3 or 8 Th isolation windows, respectively, or in
1052 triplicate 14 min or 30 min single-shot measurements from a HeLa sample acquired in DDA
1053 (left) or DIA (right) on an Orbitrap Astral²⁹. FDR was controlled at 1% at the global peptide
1054 group-level. Match between runs was used for DDA data and for DIA data, peptides were
1055 quantified irrespective of their run-specific FDR. **(B)** Peptide-level CV distributions for shared
1056 peptide groups from the same data as in **(A)**. **(C)** Peptide-level log₂-ratio density plots for
1057 shared peptide groups from the same data from the LFQbench dataset as in **(A)**. **(D)** Scatter
1058 plot of peptide-level apex quantification between CHIMERYs' MS2-based quantification (x-
1059 axis) and Minora's MS1-based quantification (y-axis) for a DIA raw file from the LFQbench
1060 dataset.

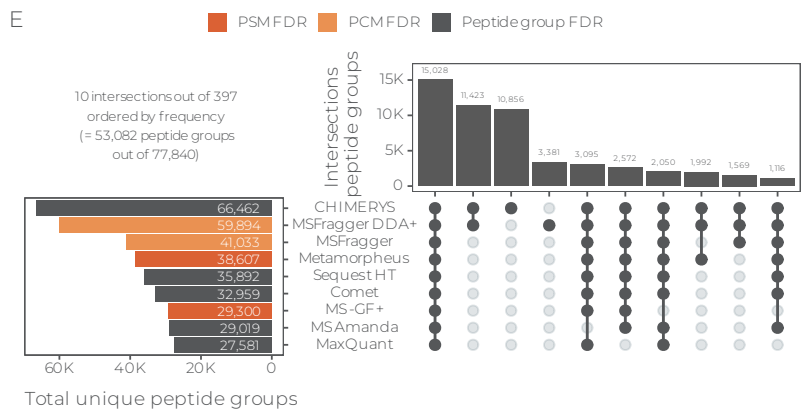
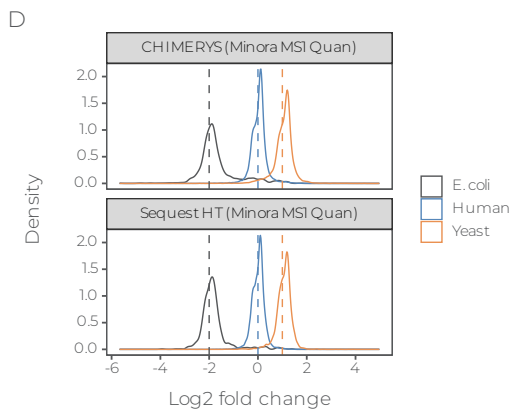
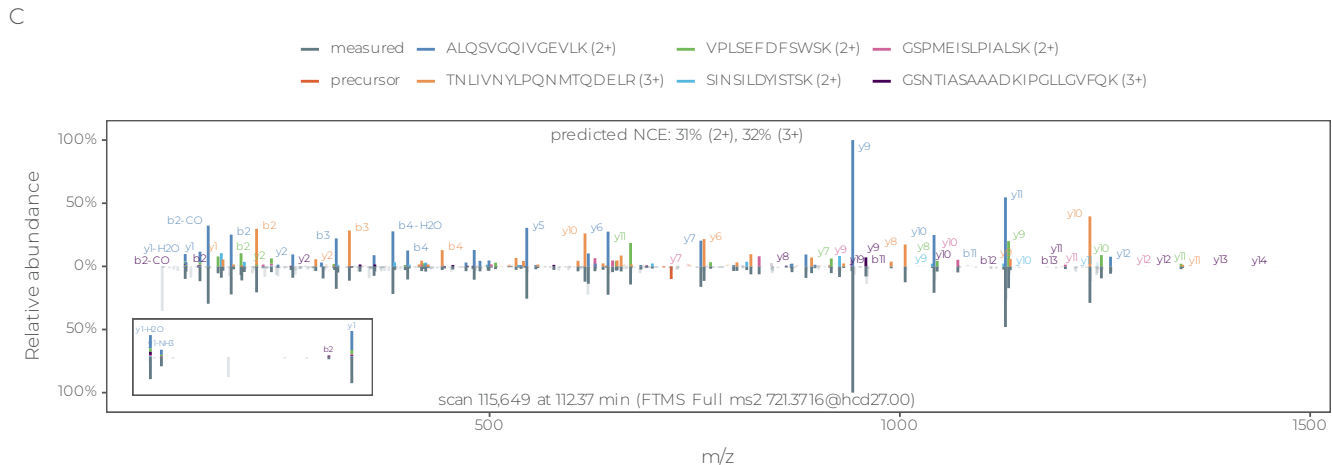
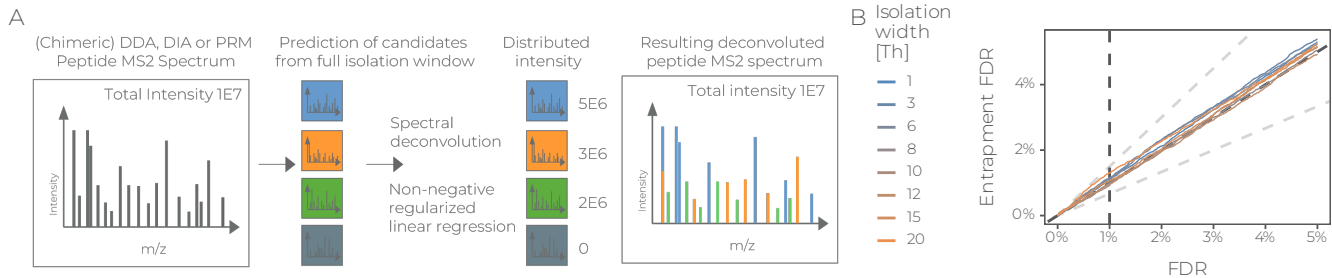
1061 References

1062

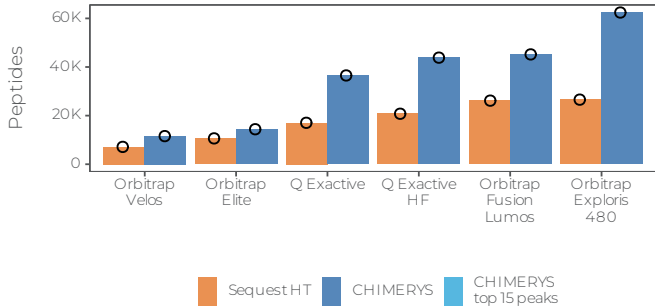
- 1063 1. Bantscheff, M., Schirle, M., Sweetman, G., Rick, J. & Kuster, B. Quantitative mass
1064 spectrometry in proteomics: a critical review. *Anal. Bioanal. Chem.* **389**, 1017–1031 (2007).
- 1065 2. Bantscheff, M., Lemeer, S., Savitski, M. M. & Kuster, B. Quantitative mass spectrometry in
1066 proteomics: critical review update from 2007 to the present. *Anal. Bioanal. Chem.* **404**, 939–
1067 965 (2012).
- 1068 3. Aebersold, R. & Mann, M. Mass spectrometry-based proteomics. *Nature* **422**, 198–207
1069 (2003).
- 1070 4. Peckner, R. *et al.* Specter: linear deconvolution for targeted analysis of data-independent
1071 acquisition mass spectrometry proteomics. *Nat. Methods* **15**, 371–378 (2018).
- 1072 5. Dorfer, V., Maltsev, S., Winkler, S. & Mechtler, K. CharmeRT: Boosting Peptide
1073 Identifications by Chimeric Spectra Identification and Retention Time Prediction. *J. Proteome*
1074 *Res.* **17**, 2581–2589 (2018).
- 1075 6. Ting, Y. S. *et al.* Peptide-Centric Proteome Analysis: An Alternative Strategy for the
1076 Analysis of Tandem Mass Spectrometry Data*. *Mol. Cell. Proteom.* **14**, 2301–2307 (2015).
- 1077 7. Fernández-Costa, C. *et al.* Impact of the Identification Strategy on the Reproducibility of
1078 the DDA and DIA Results. *J. Proteome Res.* **19**, 3153–3161 (2020).
- 1079 8. Tyanova, S., Temu, T. & Cox, J. The MaxQuant computational platform for mass
1080 spectrometry-based shotgun proteomics. *Nat. Protoc.* **11**, 2301–2319 (2016).
- 1081 9. Pino, L. K. *et al.* The Skyline ecosystem: Informatics for quantitative mass spectrometry
1082 proteomics. *Mass Spectrom. Rev.* **39**, 229–244 (2020).
- 1083 10. Tsou, C.-C. *et al.* DIA-Umpire: comprehensive computational framework for data-
1084 independent acquisition proteomics. *Nat. Methods* **12**, 258–264 (2015).
- 1085 11. The, M., Samaras, P., Kuster, B. & Wilhelm, M. Reanalysis of ProteomicsDB Using an
1086 Accurate, Sensitive, and Scalable False Discovery Rate Estimation Approach for Protein
1087 Groups. *Mol. Cell. Proteom.* **21**, 100437 (2022).
- 1088 12. Ma, K., Vitek, O. & Nesvizhskii, A. I. A statistical model-building perspective to
1089 identification of MS/MS spectra with PeptideProphet. *BMC Bioinform.* **13**, S1–S1 (2012).
- 1090 13. Fondrie, W. E. & Noble, W. S. mokapot: Fast and Flexible Semisupervised Learning for
1091 Peptide Detection. *J. Proteome Res.* **20**, 1966–1971 (2021).
- 1092 14. The, M., MacCoss, M. J., Noble, W. S. & Käll, L. Fast and Accurate Protein False
1093 Discovery Rates on Large-Scale Proteomics Data Sets with Percolator 3.0. *J. Am. Soc. Mass*
1094 *Spectrom.* **27**, 1719–1727 (2016).
- 1095 15. Gessulat, S. *et al.* Prosit: proteome-wide prediction of peptide tandem mass spectra by
1096 deep learning. *Nat. Methods* **16**, 509–518 (2019).

- 1097 16. Degroeve, S. & Martens, L. MS2PIP: a tool for MS/MS peak intensity prediction.
1098 *Bioinformatics* **29**, 3199–3203 (2013).
- 1099 17. Zhou, X.-X. *et al.* pDeep: Predicting MS/MS Spectra of Peptides with Deep Learning.
1100 *Anal. Chem.* **89**, 12690–12697 (2017).
- 1101 18. Tibshirani, R. Regression Shrinkage and Selection Via the Lasso. *J. R. Stat. Soc. Ser. B: Stat. Methodol.* **58**, 267–288 (2018).
- 1103 19. Lu, Y. Y., Bilmes, J., Rodriguez-Mias, R. A., Villén, J. & Noble, W. S. DIAMeter: matching
1104 peptides to data-independent acquisition mass spectrometry data. *Bioinformatics* **37**, i434–
1105 i442 (2021).
- 1106 20. Matzinger, M. *et al.* Micropillar arrays, wide window acquisition and AI-based data
1107 analysis improve comprehensiveness in multiple proteomic applications. *Nat. Commun.* **15**,
1108 1019 (2024).
- 1109 21. Truong, T. *et al.* Data-Dependent Acquisition with Precursor Coisolation Improves
1110 Proteome Coverage and Measurement Throughput for Label-Free Single-Cell Proteomics**. *Angew. Chem. Int. Ed.* **62**, e202303415 (2023).
- 1112 22. Puyvelde, B. V. *et al.* A comprehensive LFQ benchmark dataset on modern day
1113 acquisition strategies in proteomics. *Sci. Data* **9**, 126 (2022).
- 1114 23. Demichev, V., Messner, C. B., Vernardis, S. I., Lilley, K. S. & Ralser, M. DIA-NN: neural
1115 networks and interference correction enable deep proteome coverage in high throughput.
1116 *Nat. Methods* **17**, 41–44 (2020).
- 1117 24. Muntel, J. *et al.* Surpassing 10 000 identified and quantified proteins in a single run by
1118 optimizing current LC-MS instrumentation and data analysis strategy. *Mol. Omics* **15**, 348–
1119 360 (2019).
- 1120 25. Rosenberger, G. *et al.* Statistical control of peptide and protein error rates in large-scale
1121 targeted DIA analyses. *Nat. methods* **14**, 921–927 (2017).
- 1122 26. Baker, C. P., Bruderer, R., Abbott, J., Arthur, J. S. C. & Brenes, A. J. Optimizing
1123 Spectronaut Search Parameters to Improve Data Quality with Minimal Proteome Coverage
1124 Reductions in DIA Analyses of Heterogeneous Samples. *J. Proteome Res.* (2024)
1125 doi:10.1021/acs.jproteome.3c00671.
- 1126 27. Orsburn, B. C. Proteome Discoverer—A Community Enhanced Data Processing Suite for
1127 Protein Informatics. *Proteomes* **9**, 15 (2021).
- 1128 28. Heil, L. R. *et al.* Evaluating the Performance of the Astral Mass Analyzer for Quantitative
1129 Proteomics Using Data-Independent Acquisition. *J. Proteome Res.* **22**, 3290–3300 (2023).
- 1130 29. Guzman, U. H. *et al.* Ultra-fast label-free quantification and comprehensive proteome
1131 coverage with narrow-window data-independent acquisition. *Nat. Biotechnol.* 1–12 (2024)
1132 doi:10.1038/s41587-023-02099-7.
- 1133 30. Bouwmeester, R., Gabriels, R., Hulstaert, N., Martens, L. & Degroeve, S. DeepLC can
1134 predict retention times for peptides that carry as-yet unseen modifications. *Nat. Methods* **18**,
1135 1363–1369 (2021).

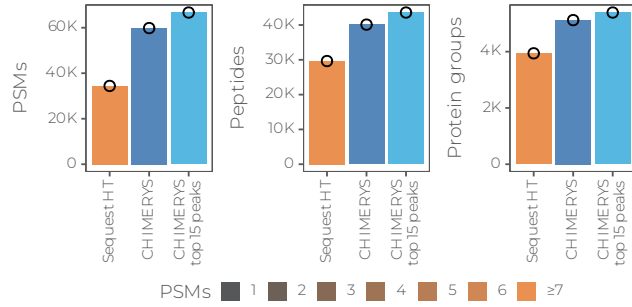
- 1136 31. Zolg, D. P. *et al.* Building ProteomeTools based on a complete synthetic human
1137 proteome. *Nat. Methods* **14**, 259–262 (2017).
- 1138 32. Vaswani, A. *et al.* Attention is All you Need. in *Advances in Neural Information*
1139 *Processing Systems 30 (NIPS 2017)* vol. 30 (Curran Associates, Inc., 2017).
- 1140 33. Chung, J., Gulcehre, C., Cho, K. & Bengio, Y. Empirical Evaluation of Gated Recurrent
1141 Neural Networks on Sequence Modeling. *arXiv* (2014) doi:10.48550/arxiv.1412.3555.
- 1142 34. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *arXiv* (2014)
1143 doi:10.48550/arxiv.1412.6980.
- 1144 35. Toprak, U. H. *et al.* Conserved Peptide Fragmentation as a Benchmarking Tool for Mass
1145 Spectrometers and a Discriminating Feature for Targeted Proteomics*. *Mol. Cell. Proteom.*
1146 **13**, 2056–2071 (2014).
- 1147 36. Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A. & Talwalkar, A. Hyperband: A Novel
1148 Bandit-Based Approach to Hyperparameter Optimization. *arXiv* (2016)
1149 doi:10.48550/arxiv.1603.06560.
- 1150 37. Zolg, D. P. *et al.* INFERYYS rescoring: Boosting peptide identifications and scoring
1151 confidence of database search results. *Rapid Commun. Mass Spectrom.* e9128 (2021)
1152 doi:10.1002/rcm.9128.
- 1153 38. Kong, A. T., Leprevost, F. V., Avtonomov, D. M., Mellacheruvu, D. & Nesvizhskii, A. I.
1154 MSFragger: ultrafast and comprehensive peptide identification in mass spectrometry-based
1155 proteomics. *Nat. Methods* **14**, 513–520 (2017).
- 1156 39. Käll, L., Storey, J. D. & Noble, W. S. qvalue: non-parametric estimation of q-values and
1157 posterior error probabilities. *Bioinformatics* **25**, 964–966 (2009).
- 1158 40. Müller, J. B. *et al.* The proteome landscape of the kingdoms of life. *Nature* **582**, 592–596
1159 (2020).
- 1160 41. Bian, Y. *et al.* Robust, reproducible and quantitative analysis of thousands of proteomes
1161 by micro-flow LC–MS/MS. *Nat. Commun.* **11**, 157 (2020).
- 1162 42. Tüshaus, J. *et al.* An optimized quantitative proteomics method establishes the cell
1163 type-resolved mouse brain secretome. *EMBO J.* **39**, e105693 (2020).
- 1164 43. Frankenfield, A. M., Ni, J., Ahmed, M. & Hao, L. Protein Contaminants Matter: Building
1165 Universal Protein Contaminant Libraries for DDA and DIA Proteomics. *J. Proteome Res.* **21**,
1166 2104–2113 (2022).
- 1167 44. Kockmann, T. & Panse, C. The rawrr R Package: Direct Access to Orbitrap Data and
1168 Beyond. *J. Proteome Res.* **20**, 2028–2034 (2021).
- 1169 45. Perez-Riverol, Y. *et al.* The PRIDE database resources in 2022: a hub for mass
1170 spectrometry-based proteomics evidences. *Nucleic Acids Res.* **50**, D543–D552 (2021).
- 1171
- 1172



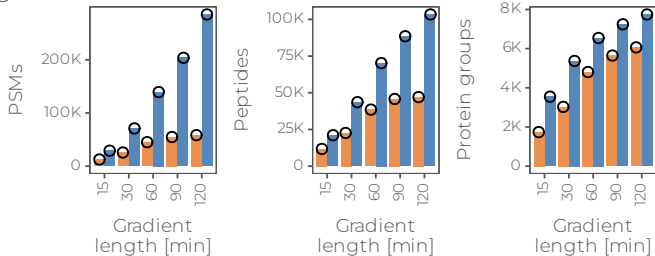
A



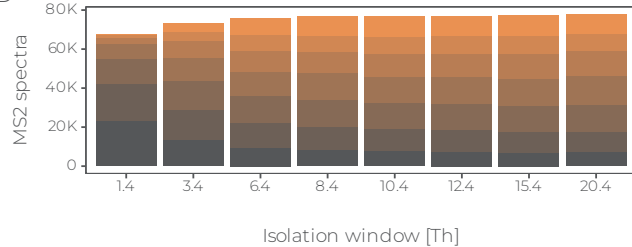
B



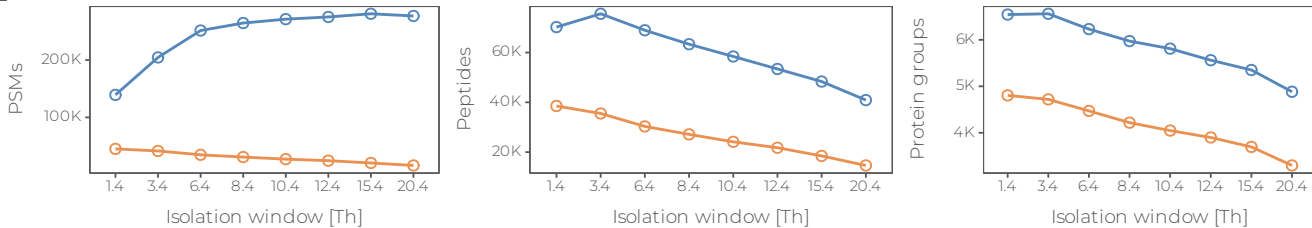
C



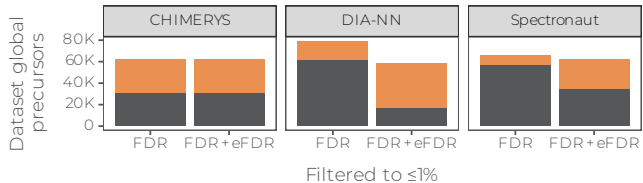
D



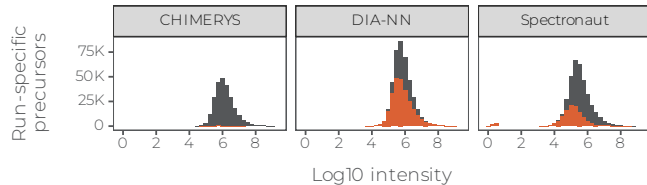
E



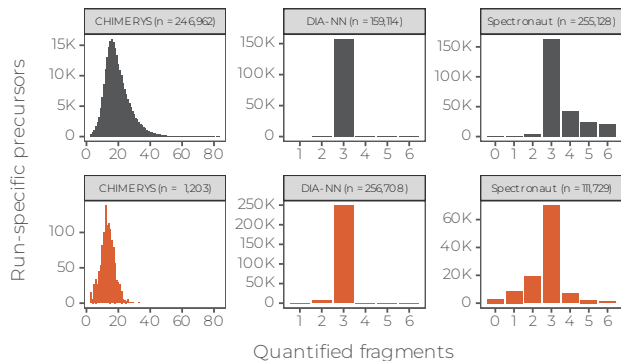
A

Quantified in ≥ 2 replicates ■ TRUE ■ FALSE

B

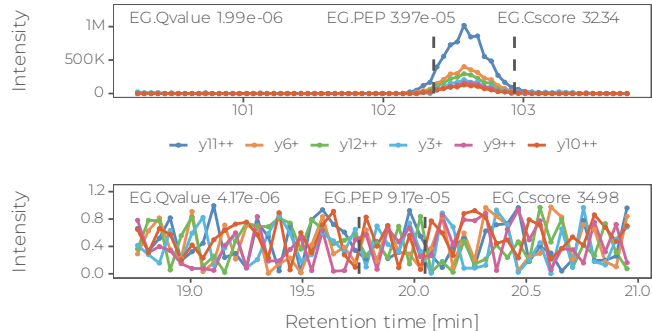
eFDR $\leq 1\%$ ■ TRUE ■ FALSE

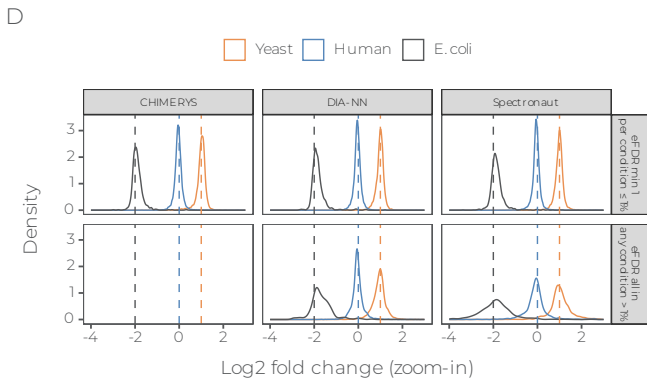
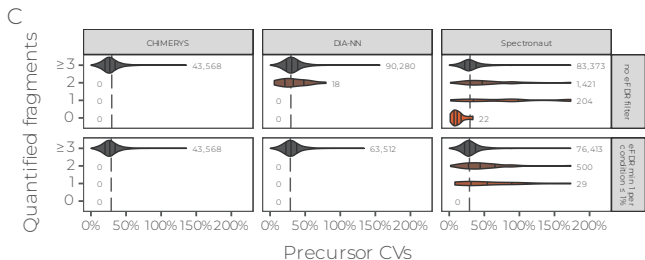
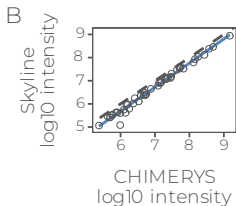
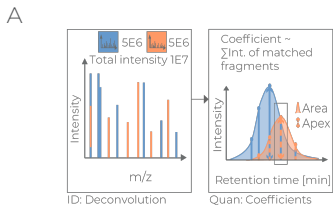
C

eFDR $\leq 1\%$ ■ TRUE ■ FALSE

D

— y5+ — y6+ — y7+ — b3+ — y8+ — y3+

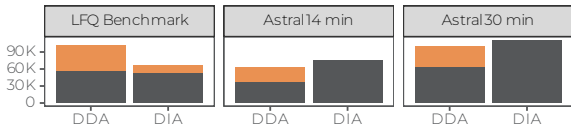




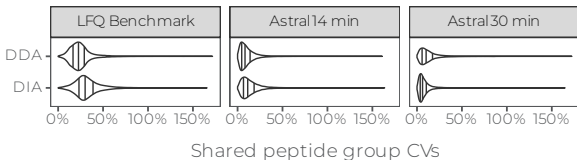
A

Quantified in ≥ 2 replicates in each condition ■ FALSE ■ TRUE

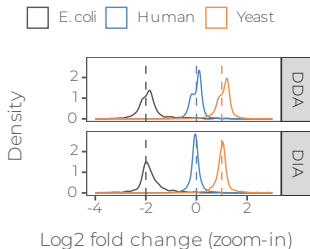
Dataset global peptide groups



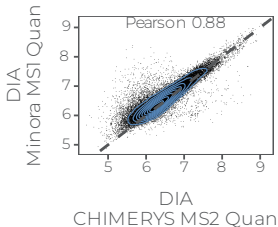
B



C



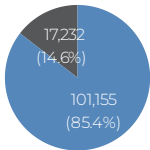
D



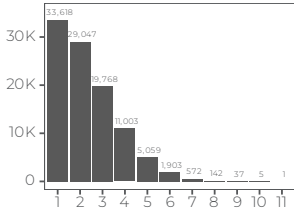
A

Identified MS2 scans

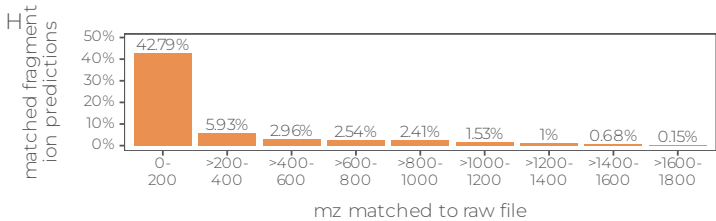
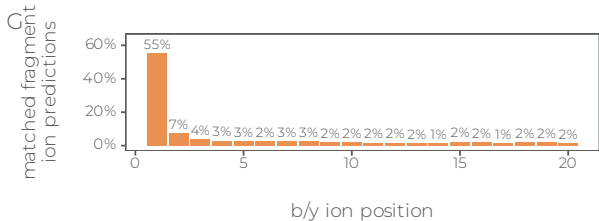
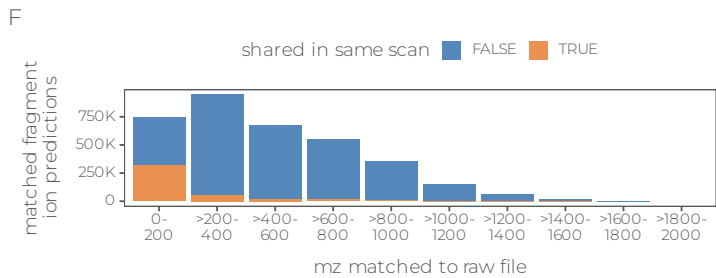
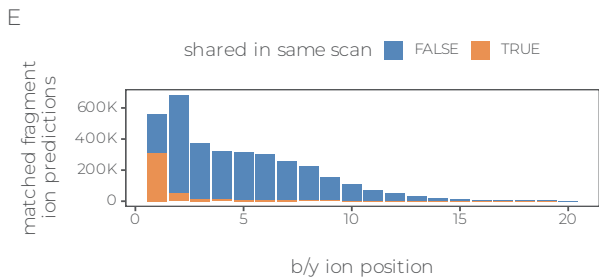
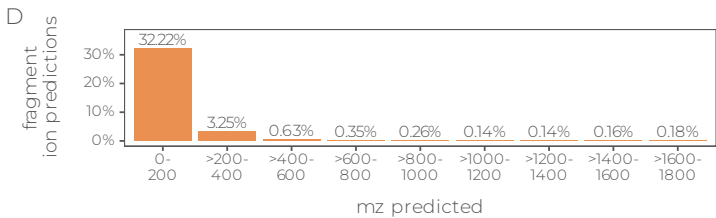
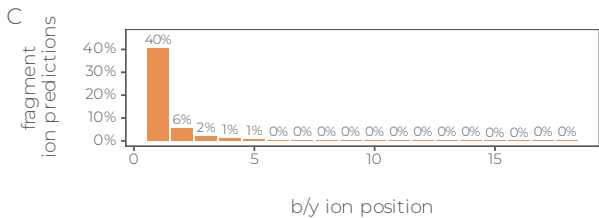
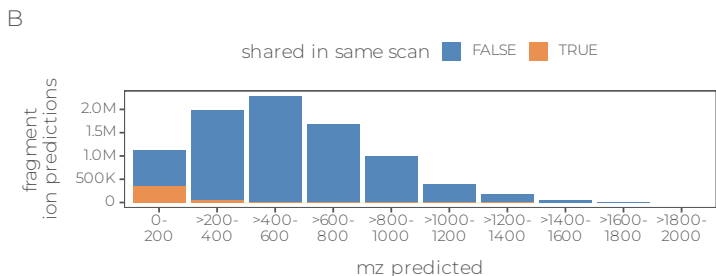
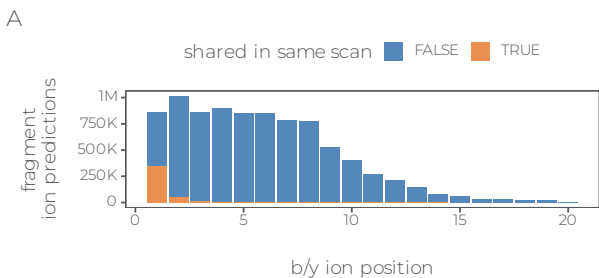
■ FALSE ■ TRUE



B

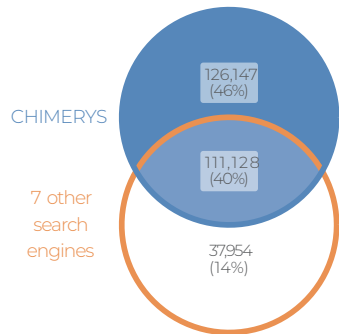
CHIMERYS
PSM count

PSMs per MS2 scan



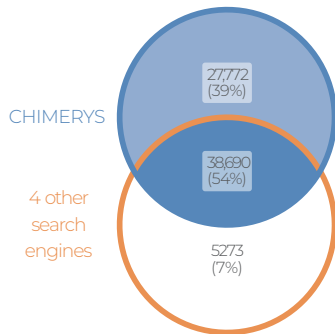
A

PSM overlap
(PSM FDR level only)



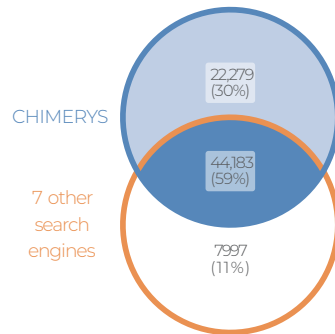
B

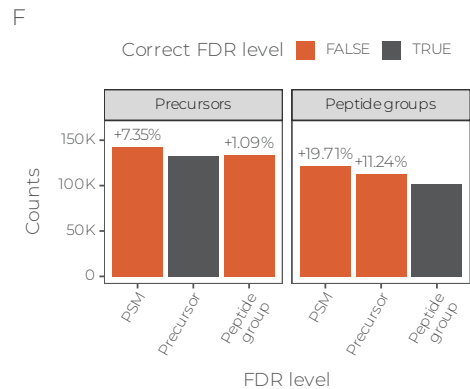
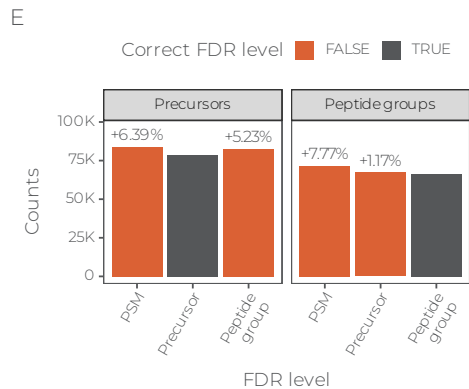
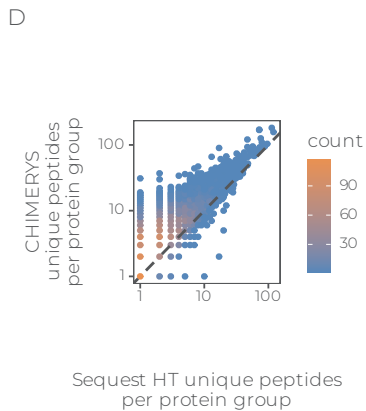
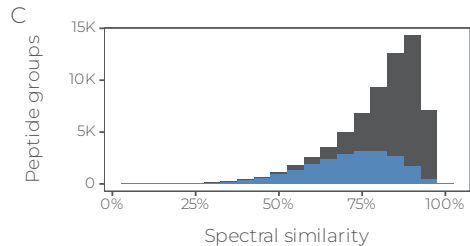
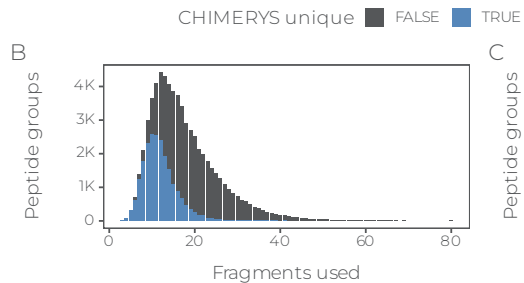
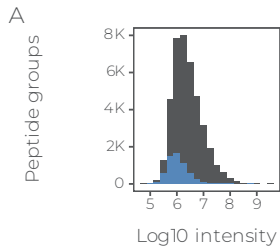
Peptide group overlap
(Peptide group FDR level only)

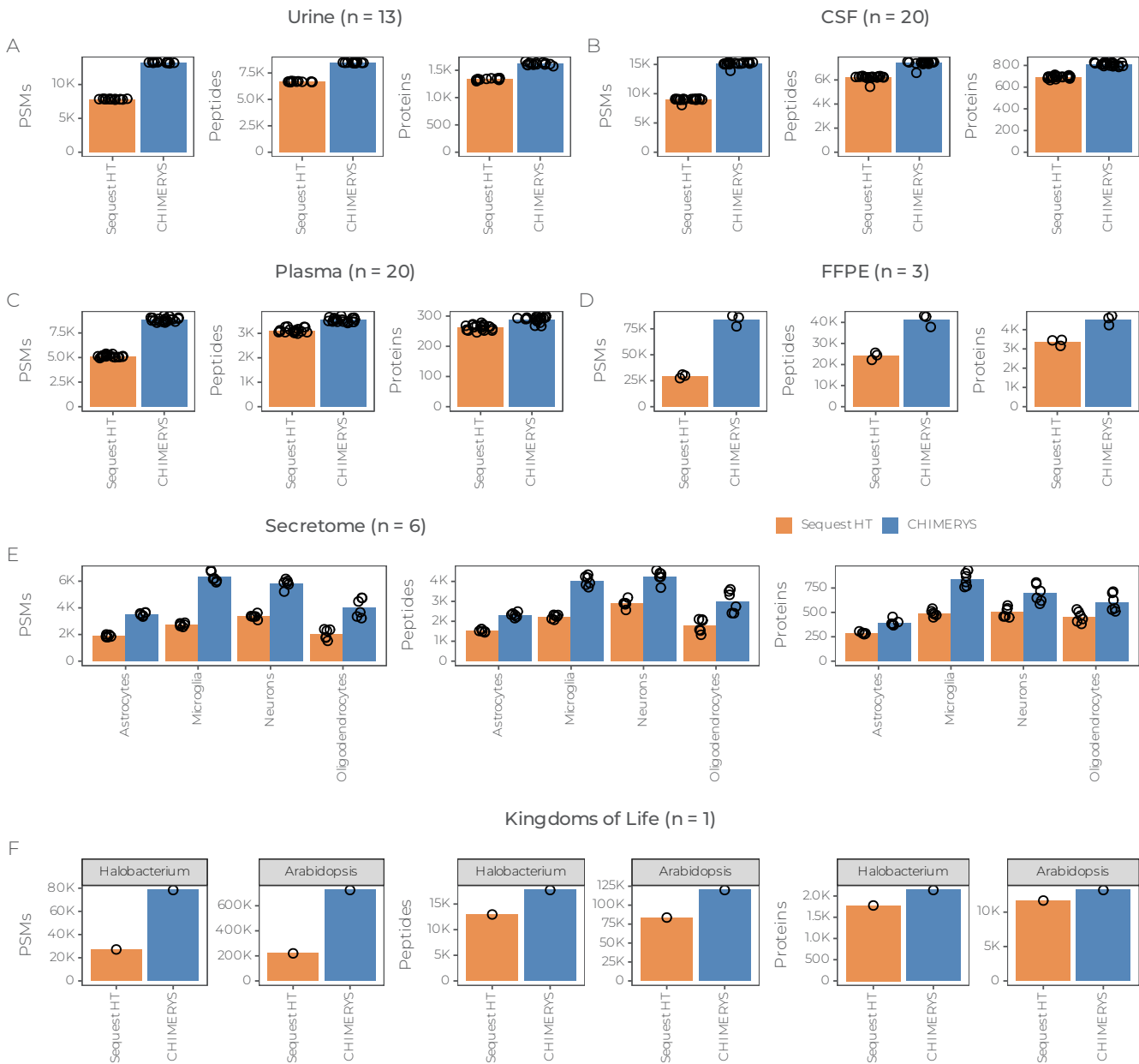


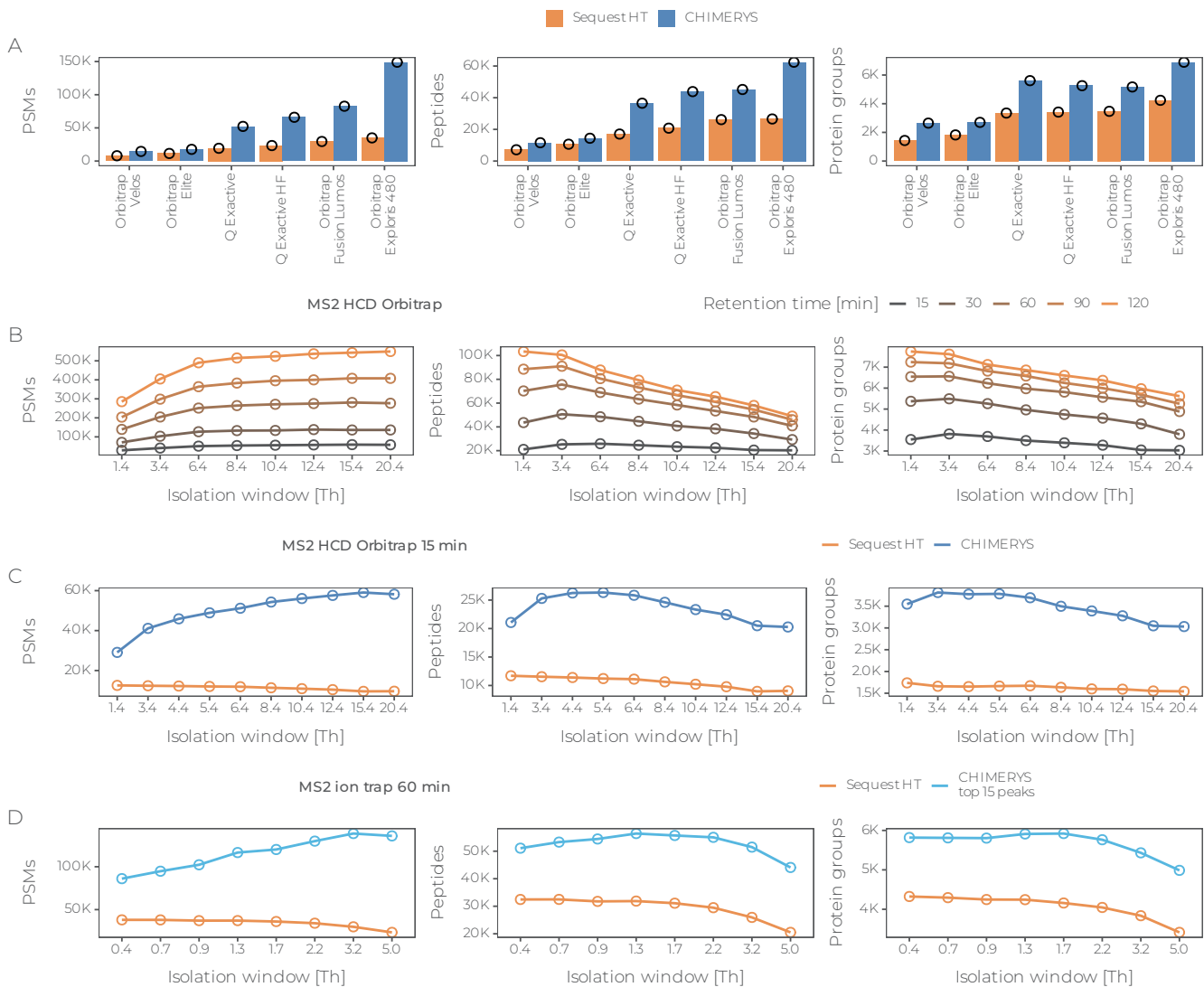
C

Peptide group overlap
(mixed FDR levels)

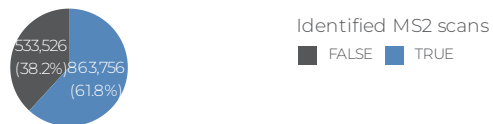




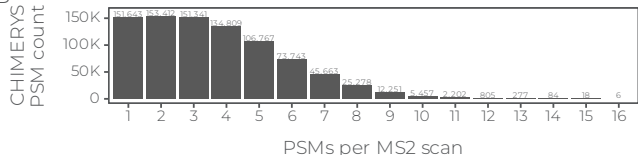




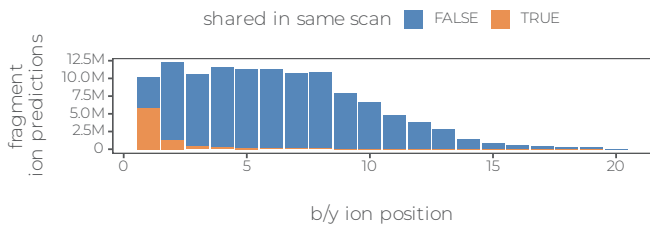
A



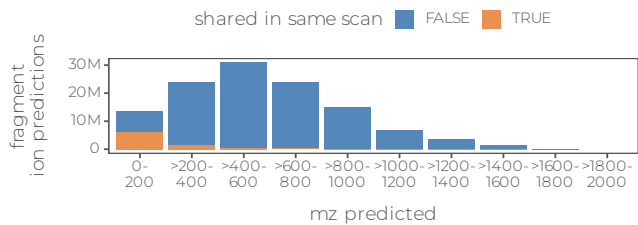
B



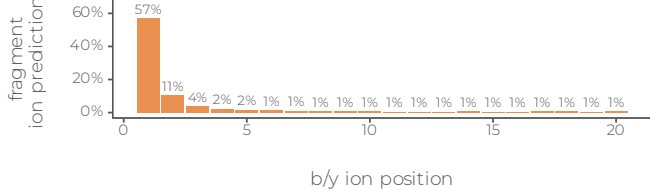
C



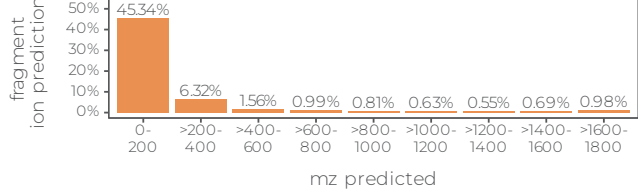
D



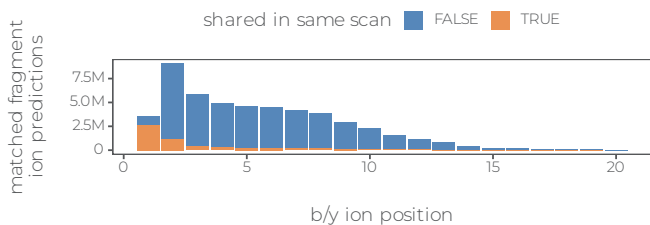
E



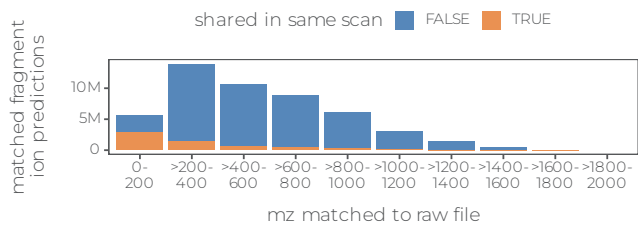
F



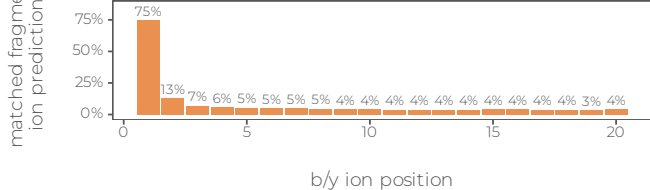
G



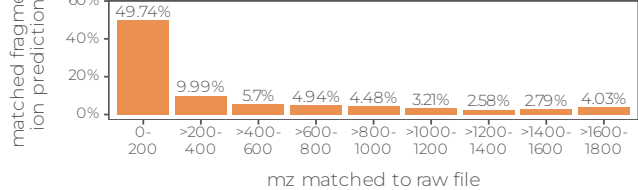
H



I

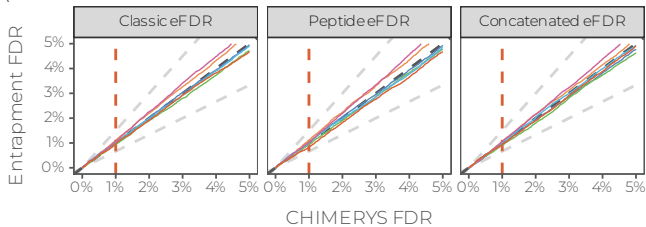


J

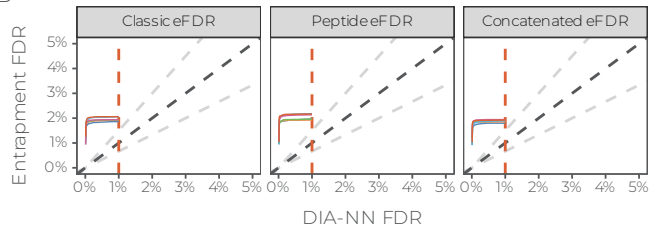


LFQ_Orbitrap_AIF_[...]

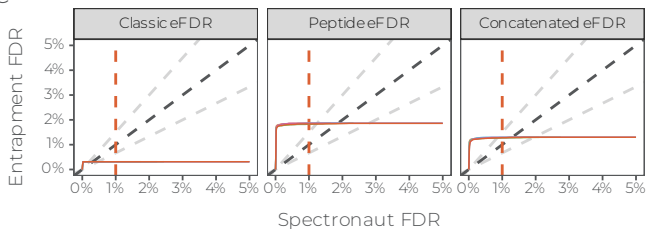
A



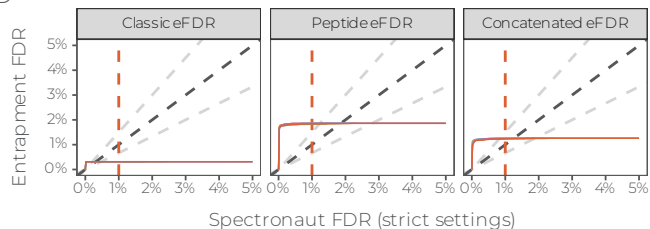
B



C

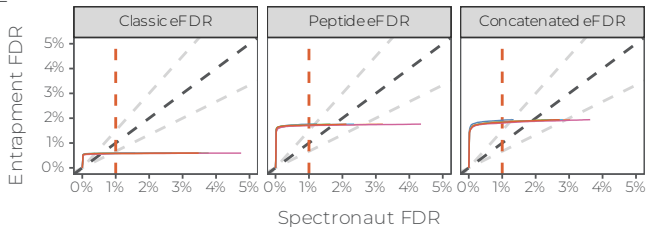


D



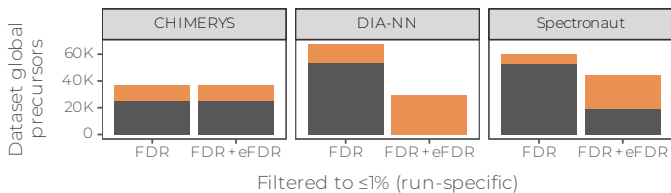
LFQ_timsTOFPro_diaPASEF_[...]

E

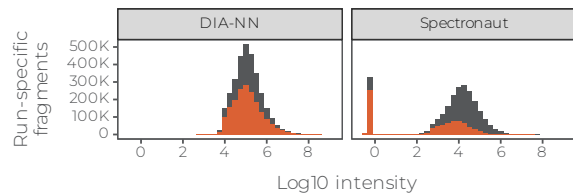


A

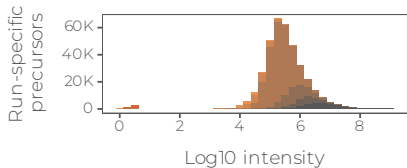
Quantified in all replicates ■ TRUE ■ FALSE



B

eFDR $\leq 1\%$ ■ TRUE ■ FALSE

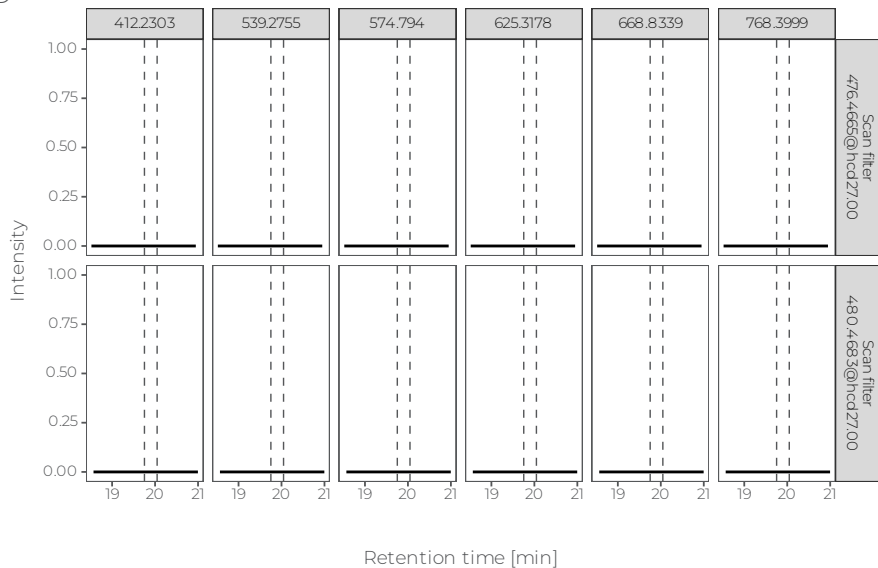
C



Curated fragments

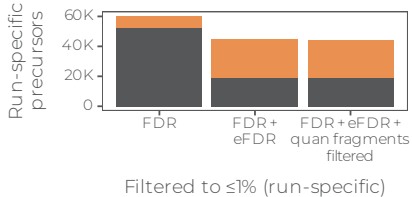


D

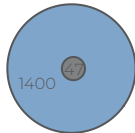


E

Quantified in all replicates ■ TRUE ■ FALSE



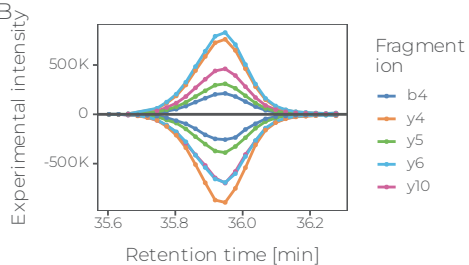
A



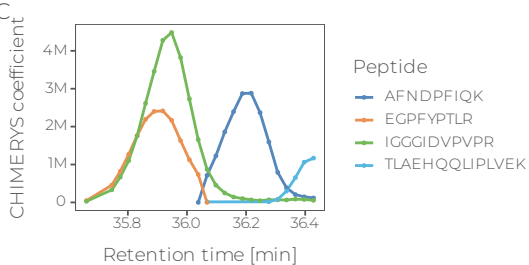
● CHIMERYs

● Skyline (manual)

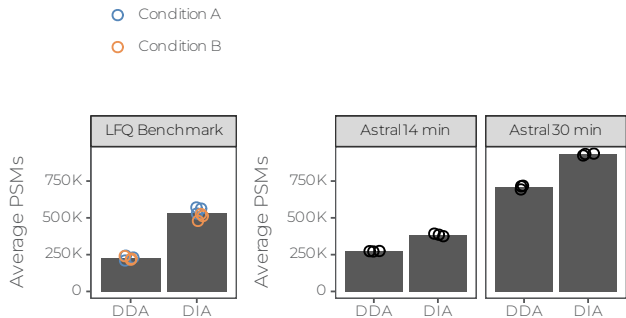
B



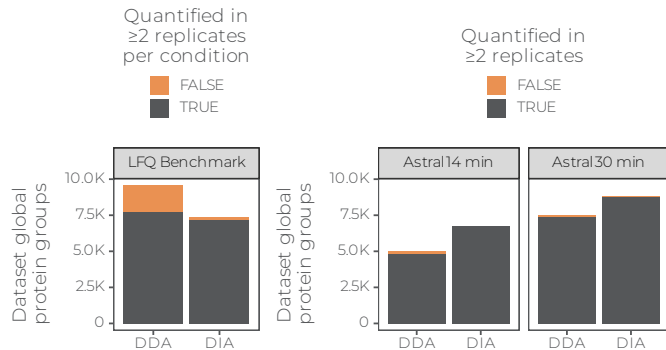
C



A



B



C

