

1 **Near-noiseless ion networks from multiple data**

2 **independent acquisition mass spectrometry**

3 **samples allow a fragment-centric perspective and**

4 **single window ion mobility DIA**

5 Sander Willems¹, Simon Daled¹, Bart van Puyvelde¹, Laura de Clerck¹,
6 Sofie vande Casteele¹, Filip van Nieuwerburgh¹, Dieter Deforce¹, and
7 Maarten Dhaenens¹

8 ¹*Laboratory of Pharmaceutical Biotechnology, Ghent University, Ghent, Belgium*

9 Poor reproducibility of data-dependent acquisition (DDA) in mass spectrometry based
10 proteomics can be complemented with data-independent acquisition (DIA). To date,
11 DIA needs to balance comprehensive data acquisition and comprehensive data analy-
12 sis. Frequently, DIA is analysed peptide-centric instead of spectrum-centric such as in
13 DDA. While peptide-centric approaches can handle chimericy and incorporate fragment
14 retention times, they unfortunately show limited false discovery rate (FDR) control.

15 Here we present single window ion mobility (SWIM)-DIA, wherein all precursors are
16 continuously fragmented in a single 2000 m/z window to acquire all fragments without
17 precursor selection or survey scans. While this provides a fully comprehensiveness data
18 acquisition, the data complexity requires deconvolution. Therefore, we leverage the
19 improved reproducibility of SWIM-DIA to convert data of multiple samples into a single
20 near-noiseless ion network based on consistent co-elution of individual ions.

21 This near-noiseless ion network offers a novel fragment-centric perspective. As a proof-
22 of-concept, we show that a naive database search algorithm can be applied to a near-
23 noiseless ion network of a mixture of commercial human, yeast and ecoli tryptic peptides.
24 Hereby, we annotated more than a 120 thousand reproducible ions as individual b - and
25 y -ions at an experimentally verified 1% FDR, accounting for 10 thousand peptides of
26 2000 proteins.

²⁷ **1 Introduction**

²⁸ Mass spectrometry (MS)-based proteomics is traditionally done with data-dependent acquisition
²⁹ (DDA). Herein MS is preceded by liquid chromatography (LC) to acquire multiple low energy
³⁰ (LE) scans for analytes that continuously elute over time. For each of these LE scans, a few
³¹ precursors are selected for high energy (HE) scans, in which fragmentation occurs. The selection
³² of precursors for fragmentation generally relies on the intensity and charge state of precursors and
³³ hence is data dependent. This acquisition methodology has several inherent limitations, such as 1)
³⁴ poor reproducibility due to stochastic precursor selection, 2) no temporal information on fragments
³⁵ obtained in HE scans, and 3) a limited duty cycle, defined as the ratio of ions formed in electrospray
³⁶ ionization (ESI) that enter the mass spectrometer and finally reach the detector, as limited time is
³⁷ spent on LE scans.

³⁸ In recent years, several data-independent acquisition (DIA) techniques have been developed that
³⁹ replace the data-dependent precursor selection of DDA with a partitioning of predefined mass over
⁴⁰ charge ratio (m/z) windows for fragmentation. The main differences between most of these tech-
⁴¹ niques are the number and size of the predefined m/z windows. Currently, the most popular
⁴² technique is probably sequential window acquisition of all theoretical mass spectra (SWATH), in-
⁴³ troduced by AB Sciex in 2012, in which each LE scan is typically followed by 32 or 64 HE scans of 20
⁴⁴ or 10 m/z wide. Another technique is Waters' MS^e that was introduced in 2004, in which each LE
⁴⁵ scan is followed by a single HE scan in which all detected precursors, typically between 0 and 2000
⁴⁶ m/z , are fragmented. They have since improved upon this technique with the introduction of an
⁴⁷ ion mobility separation (IMS) cell between their ion guide and collision cell, defining this technique
⁴⁸ as high definition MS^e (HDMS^e). The IMS cell separates precursors based on their collisional cross
⁴⁹ section (CCS). This separation, with drift time (DT) as metric, is achieved in milliseconds so that
⁵⁰ it fits exactly between the LC in which retention time (RT) is measured in seconds and the time
⁵¹ of flight (TOF) detector in which m/z is measured in microseconds. Both SWATH and HDMS^e,
⁵² among several others, have experimentally proven to provide more reproducible data than DDA
⁵³ and temporal information on fragments, while only HDMS^e has improved upon the duty cycle.

⁵⁴ Even though these seem like clear advantages of DIA compared to DDA, there is also reason
⁵⁵ for caution. Since the precursor selection is replaced with predefined m/z windows in DIA, the
⁵⁶ subsequent HE spectra are more chimeric, meaning they contain fragments from multiple precursors.
⁵⁷ This chimericity can be reduced by taking smaller windows, but this comes at a cost of needing more
⁵⁸ windows to cover fragmentation of all precursors detectable in LE. This subsequently means either
⁵⁹ shorter scan times are used for HE scans, or that the cycle time, defined as the time needed to return

60 to the same window, increases while the duty cycle decreases. The former has the disadvantage that
61 lower intensities with higher coefficient of variation (CV) are measured [1], possibly even below linear
62 detector range. The latter means fewer points can be used to define an extracted ion chromatogram
63 (XIC), typically below the recommended standard of nine points. In either case, there is a trade-off
64 between data chimericity and comprehensiveness of temporal and intensity information.

65 Once DIA data is acquired, there are several approaches to process it. Most of these approaches
66 can be divided by two characteristics: library-based versus library-free and spectrum-centric versus
67 peptide-centric. In a library-based approach a pre-annotated library with known fragmentation
68 patterns and intensities is used, as opposed to library-free approaches relying only on in-silico
69 information. Libraries introduce the potential to gain specificity from prior information, but are
70 often build with DDA data, either directly or indirectly, and thereby partially negate the advantages
71 of DIA. In spectrum-centric approaches each query spectrum is annotated with target peptides that
72 best explain these query spectra, whereas in peptide-centric approaches evidence for query peptides
73 is obtained from the acquired data. Many spectrum-centric approaches have origins in DDA and
74 are not tailor-made for DIA data. As such they can be susceptible to chimericity and therefore
75 require deconvolution to partition the spectra by precursor origin. Peptide-centric approaches on
76 the other hand tend to exhibit limited false discovery rate (FDR) control as they are often based
77 on multiple reaction monitoring (MRM)-like approaches that are not always scalable to DIA data.
78 With the exception of GROUP-DIA, all approaches are performed on a per-run basis, even though
79 the reproducibility of DIA is generally considered as its strongest trait.

80 Here, we present HistoPyA, a tool that demultiplexes DIA data into a near noiseless ion-network
81 based on replicate samples. As this ion-network has minimal noise, it eliminates the need for
82 specificity obtained through DDA-based spectral libraries. Furthermore the ion-network is neither
83 scan- nor peptide-centric, allowing a fragment-centric annotation approach developed especially for
84 DIA data with an intuitive FDR control. The creation of this ion-network builds on the following
85 hypothesis: *Poor separation of analytes is the primary cause of chimeric HE spectra. However,*
86 *there are stochastic differences between runs, even between replicate injections from a single sample*
87 *vial. As fragmentation of precursors occurs after separation (in LC and IMS or precursor window*
88 *partitioning), HE data can be demultiplexed based on (in)consistency of co-separation.* As such,
89 HistoPyA uses the greatest advantage of DIA over DDA, namely its reproducibility and temporal
90 information of fragments, as an additional dimension in the data.

91 In brief, HistoPyA consists of the following steps (Figure 1). First, raw data of each individual
92 sample is peak-picked in all dimensions to obtain a list of ions and their intensities. Each ion of
93 each sample is now defined by three coordinates: 1) an m/z apex, 2) an RT apex, and 3) a DT

94 apex or precursor window. After a quick calibration and estimation of inter-run differences, all ions
 95 of all samples are simultaneously aligned to obtain a list of aggregate ions, where each aggregate
 96 ion is composed of ions from different samples with equal m/z , RT, and DT or precursor window.
 97 Hereafter, an ion-network is created with aggregate ions as vertices and edges between aggregate
 98 ions if and only if they are consistently co-separated in all overlapping samples of their individual
 99 ions. The definition of co-separation is estimated from pairs of likely isotopes in LE scans. Finally,
 100 an X!Tandem-like hyperscore is calculated for each HE aggregate ion [2], based on the aggregate
 101 m/z of all its neighbors and itself, as well as the existence of a potential LE precursor. Using
 102 Percolator, an FDR can then be calculated for each aggregate ion, which can be extrapolated to
 103 the precursor, peptide, and protein level [3]. Optionally, a relative quantification can be performed
 104 for each aggregate ion, again allowing an extrapolation to precursor, peptide, and protein level.

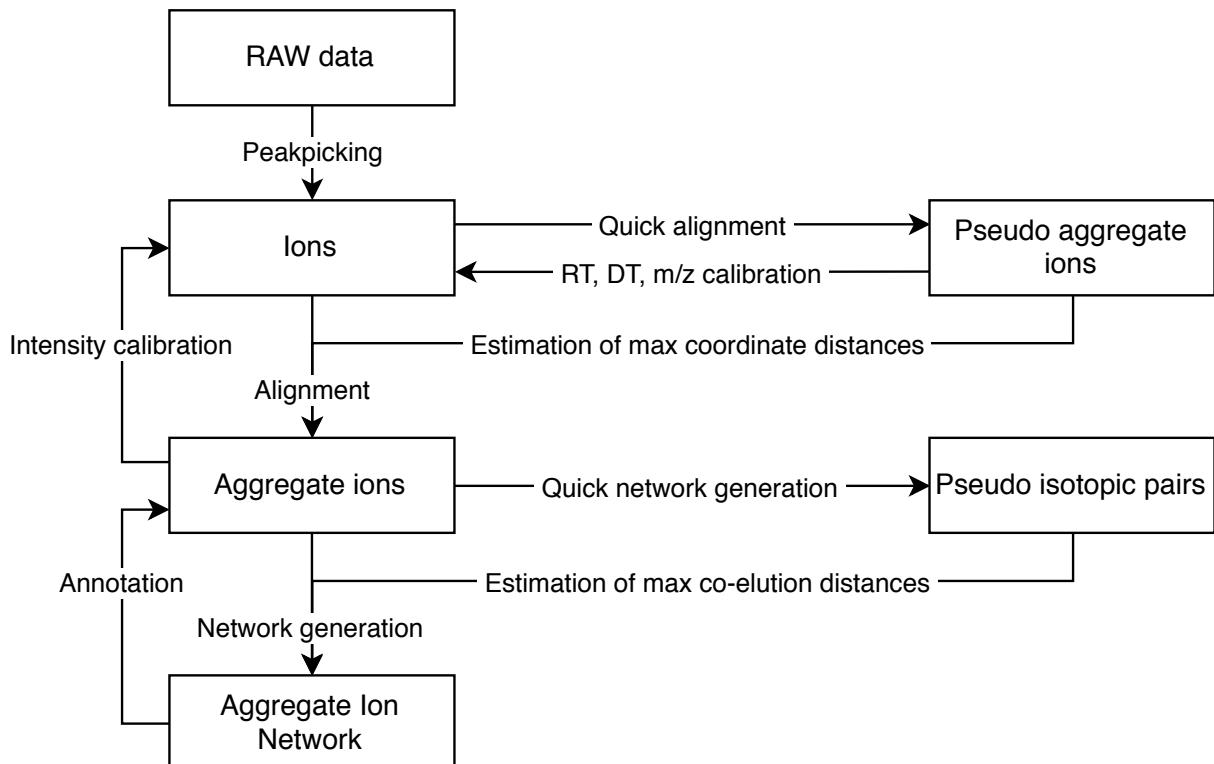


Figure 1: Schematic overview of HistoPya’s workflow.

105 To test the performance of HistoPya, two different hybrid proteomes of human, yeast and ecoli
 106 peptides were prepared with a logarithmic fold change (logFC) between condition A and B of
 107 respectively 0, 1 and -2 [4]. For each condition 3 samples were prepared, as well as a quality control
 108 (QC) sample in which all samples were pooled. For each normal sample 3 replicates were acquired
 109 and for the QC sample 9 replicates were acquired, resulting in a total of 27 samples. Based on these
 110 samples, 5 million different reproducible HE aggregate ions could be detected within a network
 111 containing 100 million edges. Depending on how well an aggregate ion was reproducible, more than

112 98% of these edges existed between ions with a similar logFC, indicating the same LE origin and
113 a very high specificity with low noise. The annotation of all these aggregate ions equally showed
114 more than 99% with a correct logFC for more than 200 thousand aggregate ions, resulting in 25
115 thousand and 4 thousand proteins, each at their respective 1% FDR.

116 These results indicate there is much to gain from such ion networks. The low noise especially
117 would allow for many new DIA annotation algorithms, including *de novo* approaches or partitioning
118 algorithms to trim the network by precursor origin to obtain traditional spectra. Moreover, the low
119 noise in the network was obtained without precursors and LE precursors are only used in the last
120 annotation step as specificity boost. This effectively means simpler hardware can be developed, in
121 which an LC is coupled directly to an IMS cell followed by a collision cell and detector, skipping any
122 precursor selection mechanisms such as a quadrupole. From a software perspective the statistics
123 could be developed further, as all aggregate ions have multiple measurement over different samples
124 and the quantification can be done on a HE level instead of a LE level. Finally HistoPyA is
125 very modular, meaning each part of the algorithm (Figure 1) can be easily substituted by other
126 approaches, allowing many improvements.

127 In DDA the measured analytes are peptides and thus a spectrum-centric or peptide-centric ap-
128 proach is justifiable. However, in DIA the direct connection between precursor- and fragment-ions is
129 not present anymore. As such, we argue that a peptide-centric approach is not fully applicable. At
130 the same time, DIA spectra are of high complexity and not independent from one another, thereby
131 also invalidating a spectrum-centric approach. What is truly measured in DIA, are independent
132 ions in multiple dimensions. As such, a more intuitive perspective than spectrum- or peptide-centric
133 would be fragment/ion-centric.

134 2 Results

135 2.1 Peak picking and inter-run ion alignment

136 The core premise of HistoPyA is that HE ions which are consistently co-eluting in all their rep-
137 resentative samples are likely to be derived from the same LE precursor. Therefore, a primary
138 component of HistoPyA is its ability to determine which ions are *the same* in different samples.

139 Using Water's Apex 3D peak picking algorithm the raw data from 27 samples (nine per condition
140 A, B and QC) were peak picked in all of its dimensions: m/z , RT, DT and intensity. Even though
141 these samples had different ratios for human, yeast and ecoli proteins, the resulting number of

142 detectable ions was similar for all of them, with each sample resulting in roughly 6 million ions in
143 both LE and HE ions (Supplementary Table 1) at the lowest thresholds of 1 for both LE and HE.

144 **2.1.1 Calibration and estimation of mass over charge ratio, retention time and drift time with**
145 **pseudo aggregate ions**

146 To compensate for small differences in acquisition, all 27 samples were calibrated with a quick
147 alignment. To perform this quick alignment the peak picked ions of all samples were merged into a
148 single list and only those ions with an intensity larger than 2^{14} were used. Of all these ions more
149 than 10 thousand pseudo aggregate ion could be formed as they had a unique m/z which was found
150 exactly once in each sample (Supplementary Figure 7). A thousand of these pseudo aggregate ions
151 had either RT or DT outliers and were removed, and the remaining pseudo aggregate ions
152 were equally partitioned in a group for calibration and a group for estimation. With the pseudo
153 aggregate ions for calibration both the m/z and DT of all peak-picked ions were corrected by at
154 most 10 parts per million (ppm) and 10000 ppm respectively, depending on their sample origins.
155 For the RT calibration the pseudo aggregate ions were grouped in more than 500 groups so that
156 each group had a distinct RT and the RT of each ion of each pseudo aggregate ion in the group was
157 strictly smaller than the RT of each ion in the next pseudo aggregate ion group. With these pseudo
158 aggregate ions groups a piece-wise linear transformation was performed to calibrate the ions of each
159 sample, resulting in a more consistent retention time (Supplementary Figure 8).

160 After this calibration, the other half of the pseudo aggregate ions were used to estimate the
161 maximum inter-run distance for each of RT, m/z , DT, resulting in respectively 5 ppm, 6000 ppm and
162 0.2 minutes (Supplementary Figures 9, 10 and 11). These estimates are relatively small, indicating
163 that the calibration performed well or that the acquisition was very robust. While the latter is
164 subjective to evaluate without comparison, the estimates with uncalibrated coordinates are more
165 than twice as large (Supplementary Figures 12, 13 and 14), indicating that at least the former
166 statement is true.

167 **2.1.2 Inter-run ion alignment**

168 With the estimation parameters obtained from the pseudo aggregate ions, more than 100 million
169 pairs of ion neighbors could be defined. However, the definition of neighboring implies a transitive
170 relation over a path of multiple neighboring pairs. With this transitivity, there are clusters with over
171 a 100 ions in them, sometimes with a sample being represented 10 times. To avoid such ambiguous

172 situations, the pairs of ion neighbors were trimmed so that each remaining cluster, hereafter defined
173 as an aggregate ion, has at most one ion per sample. With this trimming 10 million pairs of
174 ion neighbors were removed. This relatively low percentage of trimming implies the estimation
175 parameters are stringent enough to obtain a good specificity. Most of the resulting aggregate ions
176 were either fully reproducible or random noise (Supplementary Figures 15 and 16). As expected,
177 the fully reproducible aggregate ions have the highest average intensities, while the noisy aggregate
178 ions are less intense (Supplementary Figures 17 and 17).

179 **2.1.3 Intensity calibration and validation**

180 Intensity was not used to create the aggregate ions and was therefore not calibrated with the pseudo
181 aggregate ions. This calibration is done after the full alignment to include as much information as
182 possible. On average, this calibration reduced the intensity CV of anchors by a factor 3 (Supple-
183 mentary Figures 19, 20, 21, 22, 23 and 24), resulting in e.g. more than 80% of all fully reproducible
184 anchors with a CV below 20 for each of condition A, B and QC.

185 With this calibrated intensity, the logFC values between condition A and B could be determined.
186 As expected, these logFC values indeed show three groups of -2, 0 and 1, corresponding to the
187 original mix ratios of the samples (Supplementary Figure 25).

188 **2.1.4 Robustness to noise**

189 To assess HistoPyA's ability to deal with noise, different peak picking thresholds in Waters' Apex3D
190 software were tested. For nearly all partially reproducible aggregate ions, additional samples could
191 be found at lower ion count thresholds. Fully reproducible aggregate ions were not hindered by the
192 potential extra interference (Supplementary figure 26). However, halving the ion count threshold in
193 Apex3D results in doubling the amount of noisy ions in HistoPyA. Nearly all ions are reproducible
194 signal at an ion count threshold of 100.

195 **2.2 Aggregate ion network**

196 With the 5 million aggregate ions obtained at a peak picking ion count threshold of 1, a quick
197 assessment of potential isotopic pairs was made. From the 1 million fully reproducible aggregate
198 HE ions, 100 thousand pairs of pseudo isotopic pairs were obtained. Based on these isotopes, the
199 maximum difference in RT and DT peaks was estimated to be smaller than 0.05 minutes and 7000

200 ppm respectively for all of the samples (Supplementary Figures 27 and 28).

201 With these maximum difference per sample the potential *consistent* co-elution between all HE
202 aggregate ions was determined. Hereby, a total 400 million pairs of aggregate ions could be defined
203 as *consistently* co-eluting, implying a probable same precursor in LE. While there are some aggregate
204 ions with more than a 100 neighbors, most aggregate ions have between 1 and 13 neighbors showing
205 a high specificity (Supplementary Figure 29). Of note, 15% of all anchor ions do not have a single
206 neighbor and could be considered noise, even though they are reproducible over multiple runs.

207 Based on the logFC values between condition A and B many aggregate ions can be classified
208 by most likely organism origin as human, yeast or ecoli. Since neighboring aggregate ions are
209 expected to come from the same precursor, they are also expected to have the same organism
210 origin. Moreover, neighboring aggregate ions with a different logFC are most likely false positives.
211 While neighboring aggregate ions with a small sample overlap have many false positives, this false
212 positive rate decreases when aggregate ions overlap in more samples. Especially fully reproducible
213 aggregate ions that are neighbors show a low false positive rate of less than 1% (Figure 2).

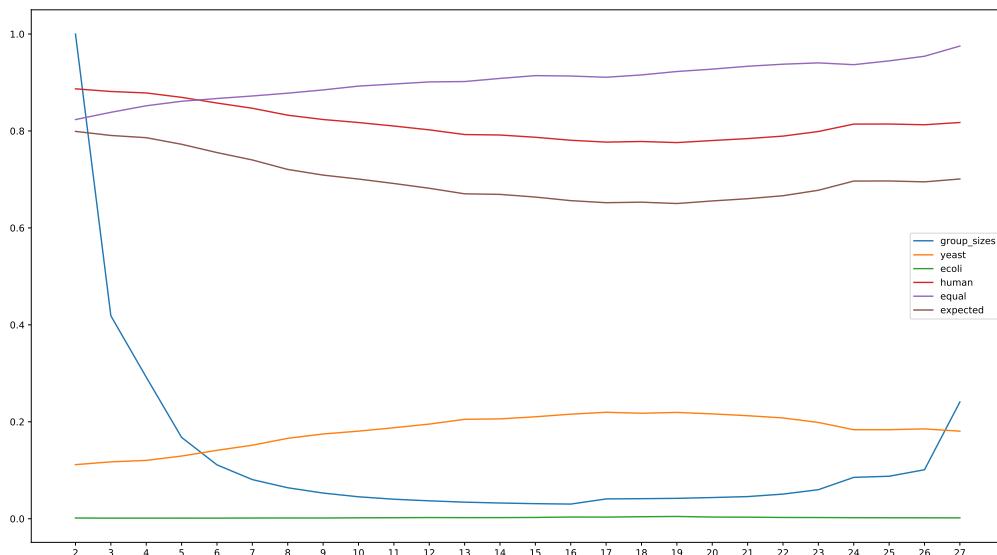


Figure 2: Neighbor organism origins.

214 Add a chimeric network (@m/z) that declusters by reproducibility. When HE aggregate ions are
215 deisotoped based on their neighbors, the delta mass frequencies are most prevalent for amino acid
216 delta masses (Supplementary methods TODO). When only LE aggregate ions are considered and
217 deisotoped, the singly charged aggregate ions have drift time that confirm the charge state as being
218 1.

219 Delta masses are amino acids, ergo the neighbors belong together.
220 To determine consistent co-elution between LE and HE, only consistency in RT was considered.
221 This was essential because there is a DT shift between LE and HE scans of roughly 3 (Supplementary
222 Figure 31)

223 TODO

224 **2.3 Fragment-centric Annotation**

225 Based on the ion network, a ion count (noiseless) was calculated for each individual ion. For each
226 possible b- or y-ion explanation of a single fragment, the count of neighboring ions from the same
227 peptide was determined. The frequency of all unique counts was determined and log transformed,
228 and the e-value for the most extreme count was determined similar as to XTandem. No intensity
229 was used (ion network). As opposed to spectrum-centric, i.e. precursor-centric. We do not use
230 presurors, but fragments create similar specificity... (PNM as opposed to PIMs) This leaves us with
231 PTMs and SNPs, i.e. open mass search. Which we will address through m/z estimation with dt.

232 **2.3.1 Model validation**

233 Count distributions (NOVEL) follow good linear distribution with outlier (XTandem).

234 Fragment-centric annotation of DDA data gives similar results as Mascot annotation. Precusor-
235 less?

236 **2.3.2 Results**

237 96000 PNM on 95000 fragments were significant out of 10e06. 1+ is annotated. But reproducibility
238 filtering: 100000 on 1e6 (plot). More DIA peptides found than DDA. All highest scoring PNM are
239 retained per fragment ion (identical number of ion counts matched peptide)

240 Dechimerization (4).

241 We annotate 10% (40.000 are annotated)

242 Peptide annotation goes through percolator, i.e. best scoring PNM = peptide score Overview of
243 fragment features (used for percolator). Percolator increased significance by many percent.

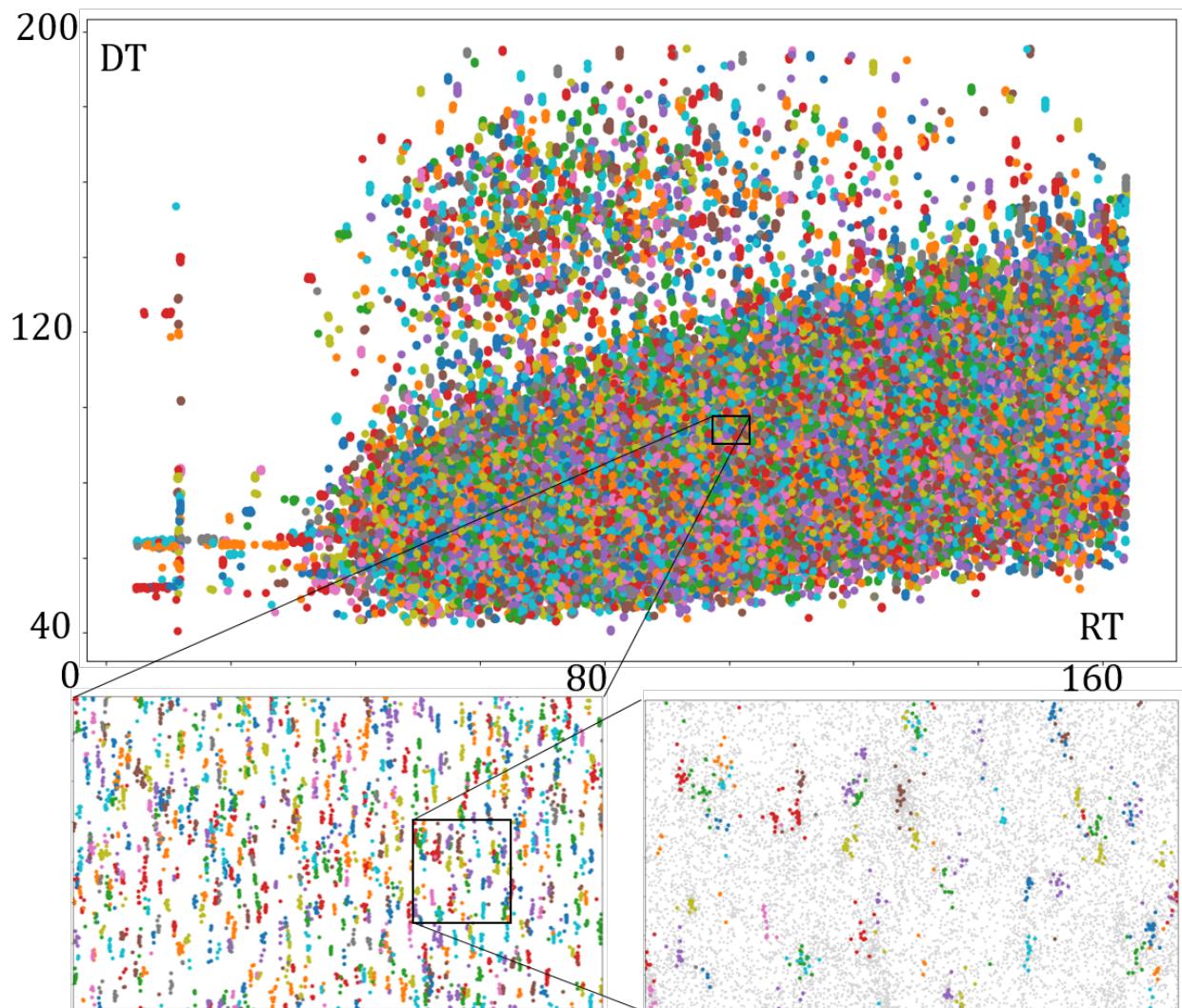


Figure 3: Annotated fragments.

244 Precursor existence filter very useful, but not essential. Alternatively You percolate your precursor
245 mass out of drift (allowing for SWIM)

246 Unannotated stuff can be annotated thanks to noiselessness, i.e. fully reproducibles (400.000)

247 250.000 have an annotated neighbor = unannotated ions: isotopes: 20.000 => 15% Non-b-y-ions
248 neutral losses non-significant PNMs

249 MEANING 150.000 belong to unannotated clusters = unannotated peptides (ISD, PTMs, missed
250 cleavage,...):

251 AND THERE ARE 41 loners

252 Vendor specs of >12000 peptides for human reached.

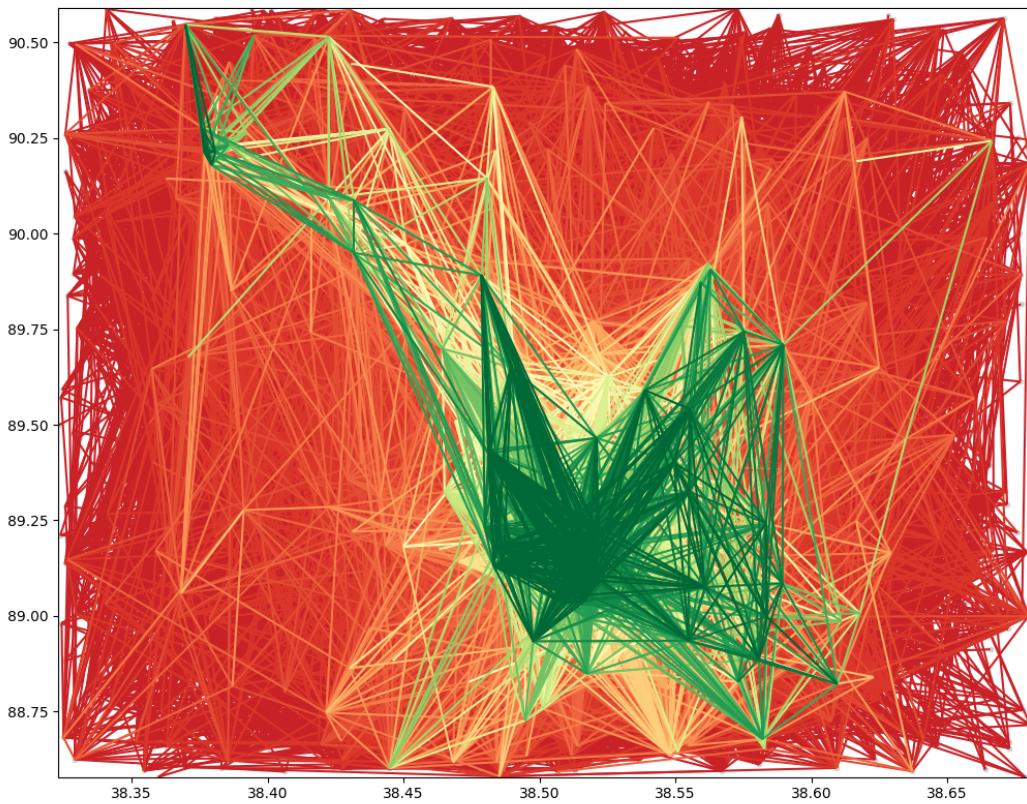


Figure 4: Dechimerization.

253 **2.3.3 Annotation validation**

254 Annotated (significant and by proxy) fragment/precursor/peptide organisms are correct with LFQ
 255 organism classification.

256 RTs coincide with mascot DDA rts.

257 Decoy-decoy (search against pyrococcus) gives no hits.

258 **2.4 Fragment-centric Quantification**

259 You have the experiment in a single network: performance Specificity with reduced interference?
 260 Removing stochasticity give ultra clean data Intensity is for Quant, not for any other purpose

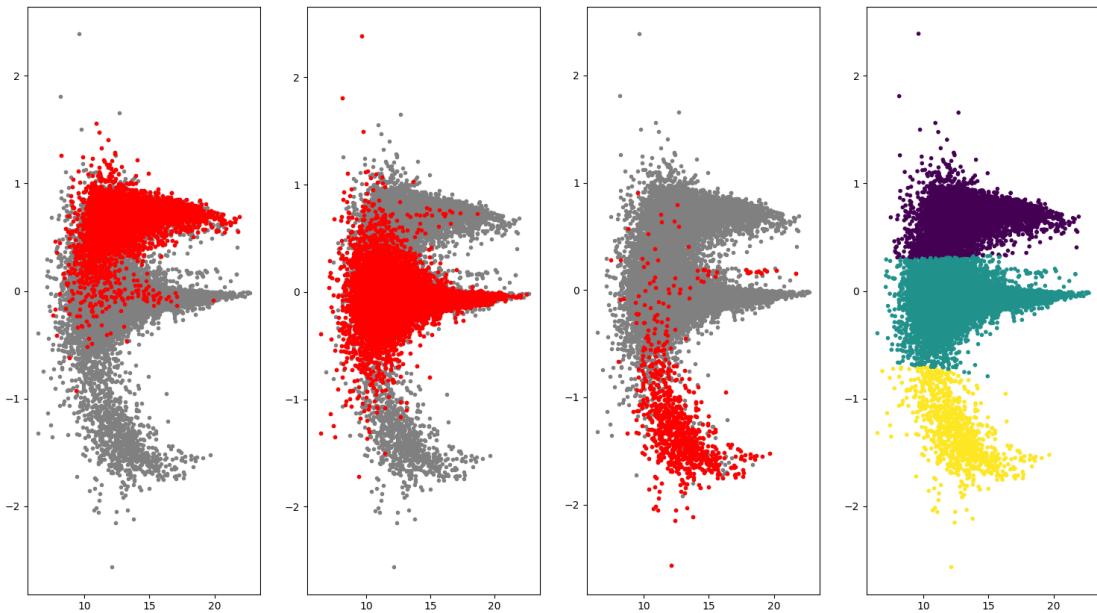


Figure 5: Experimental annotation and theoretical organism classification

261 **2.5 Fragment-centric Application**

262 Public data? LOPIT Kathryn Lilley

263 **2.6 SWIM**

264 Precursor existence filter very useful, but not essential. Alternatively You percolate your precursor
265 mass out of drift (allowing for SWIM)

266 **3 Discussion**

267 Creating pseudo MSMS spectra for traditional annotation tools (without precursor)?

268 HistoPyA is modular, meaning peak picking, calibration, ion alignment, aggregate ion network
269 creation, annotation van all be easily replaced.

270 Cosmology vs quantum mechanics

271 Noiseless data simplifies many annotation algorithms, including de novo for DIA.

272 In theory, many replicates (with high overlap requirements) would partition the aggregate ion
273 network in fragment groups all belonging to the same precursor (noise-free pseudo MSMS spectra).

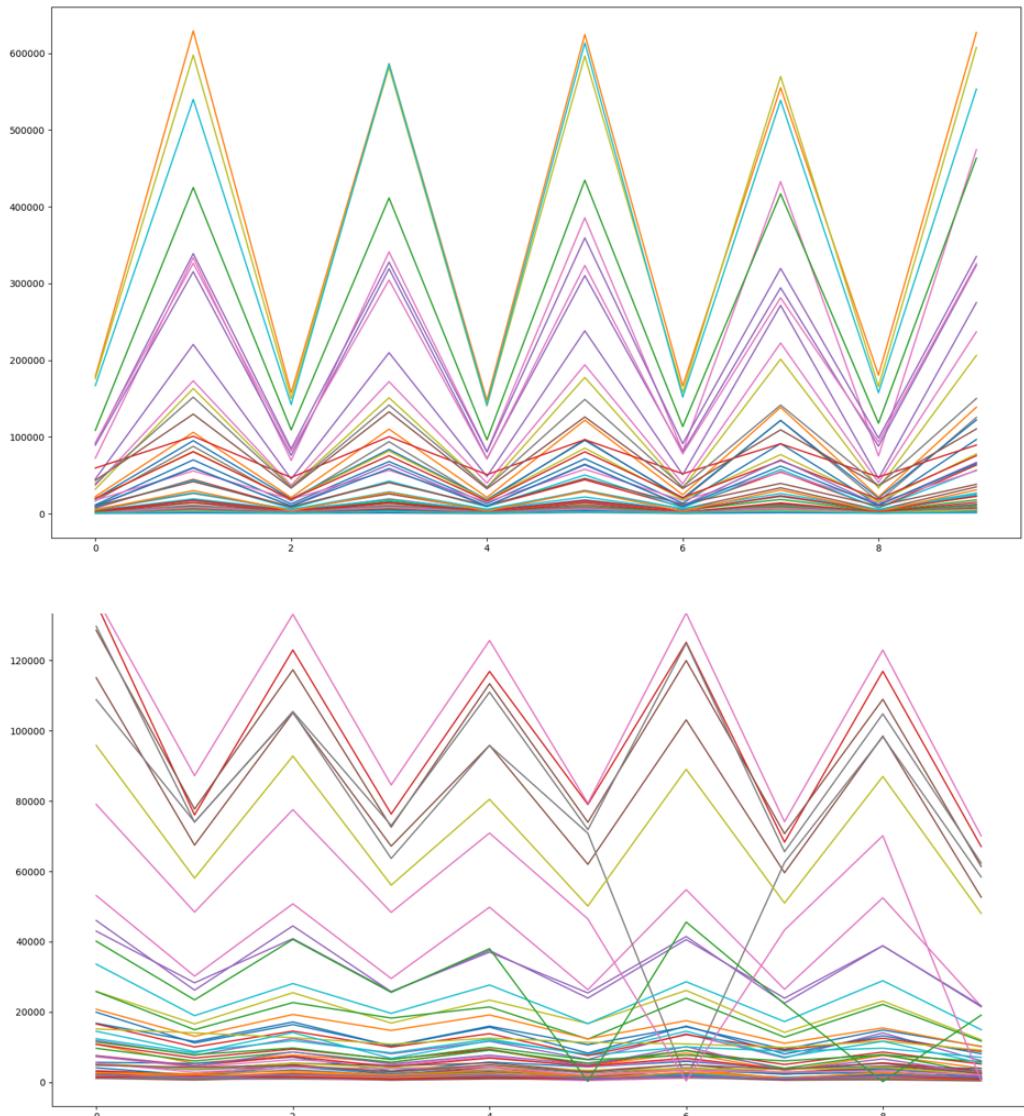


Figure 6: fragment quantifications per peptide

274 Results are better than DDA and more transparent (especially FDR) than other DIA approaches.

275 HE only could allow better/simpler/cheaper instrumentation.

276 Annotations match MS2PIP theoreticals? We start from aggregate spectra, as is done in other

277 tools as well

278 **Data and software availability**

279 All data is available at the ProteomeXchange consortium with identifier TODO. This includes raw
280 data and data peakpicked with Waters' commercial Apex3D software. Complete algorithmic results
281 as presented in this manuscript, including parameters, logs, figures, and other in/output files are
282 deposited alongside this data.

283 The complete source code for version TODO of single window ion mobility (SWIM)-DIA is
284 available at GitHub TODO. In-house scripts performing label-free quantification (LFQ) validation
285 and recreating figures are included in an additional sub folder in the GitHub repository, but require
286 original result files to be downloaded from ProteomeXchange. A minor test case illustrating how to
287 use the software on novel samples provided by the user is included in the GitHub repository.

288 QC files monitoring general MS performance are available at the Panorama website with identifier
289 TODO.

290 **Conflict of interest**

291 The authors declare no conflict of interest.

292 **Acknowledgements**

293 This research was funded by Research Foundation Flanders (FWO) research project grant G013916N,
294 FWO mandate 12E9716N (MD), FWO mandate 3F016517 (BVP), and Flanders Innovation & En-
295 trepreneurship (VLAIO) mandate SB-141209 (LDC).

296 The authors would like to express their gratitude to the Waters informatics team (including, but
297 not limited to; Hans Vissers, Scott Geromanos and Steve Cievarini) and Lennart Martens (Ghent
298 University (UGhent)) for their critical feedback throughout the project. Samples were acquired at
299 the ProGenTomics facility. Computational assistance was provided by Yannick Gansemans (UGh-
300 ent) and Laurentijn Tilleman (UGhent).

301 **Author contributions**

302 SW and MD conceived the ideas of creating noiseless ion-networks with replicates and annotating
303 this fragment-centric. SD and BVP performed all sample preparation and data acquisition. SW
304 performed all computational analysis. SW and MD wrote the draft manuscript. MD and DD
305 supervised the project. All authors provided critical feedback on the manuscript and approved the
306 final version.

307 **Acronyms**

308 **m/z** mass over charge ratio

309 **CCS** collisional cross section

310 **CPU** central processing unit

311 **CSV** comma separated values

312 **CV** coefficient of variation

313 **DDA** data-dependent acquisition

314 **DIA** data-independent acquisition

315 **DT** drift time

316 **DTT** dithiothreitol

317 **ESI** electrospray ionization

318 **FDR** false discovery rate

319 **FWO** Research Foundation Flanders

320 **GB** gigabytes TODO bits?

321 **HDMS^e** high definition MS^e

322 **HE** high energy

323 **IMS** ion mobility separation

324 **LC** liquid chromatography

325 **LE** low energy

326 **LFQ** label-free quantification

327 **logFC** logarithmic fold change

328 **MRM** multiple reaction monitoring

329 **MS** mass spectrometry

330 **PIM** peptide-ion match

331 **ppm** parts per million

332 **PSM** peptide-spectrum match

333 **ptp** point-to-point

334 **QC** quality control

335 **RAM** random-access memory

336 **RANSAC** random sample consensus

337 **RT** retention time

338 **SWATH** sequential window acquisition of all theoretical mass spectra

339 **SWIM** single window ion mobility

340 **TOF** time of flight

341 **UGhent** Ghent University

342 **VLAIO** Flanders Innovation & Entrepreneurship

343 **w/w** weight for weight

344 **XIC** extracted ion chromatogram

345 **References**

- 346 [1] Limonier F, Willems S, Waeterloos G, Sneyers M, Dhaenens M, Deforce D. Estimating the
347 reliability of low-abundant signals and limited replicate measurements through MS2 peak area
348 in SWATH. PROTEOMICS;0.ja:1800186.
- 349 [2] Craig R, Beavis RC. TANDEM: matching proteins with tandem mass spectra. Bioinformatics
350 2004;20.9:1466–1467. eprint: /oup/backfile/content_public/journal/bioinformatics/
351 20/9/10.1093/bioinformatics/bth092/2/bth092.pdf.
- 352 [3] The M, MacCoss MJ, Noble WS, Käll L. Fast and Accurate Protein False Discovery Rates on
353 Large-Scale Proteomics Data Sets with Percolator 3.0. Journal of The American Society for
354 Mass Spectrometry 2016;27.11:1719–1727.
- 355 [4] Navarro P, Kuharev J, Gillet LC, Bernhardt OM, MacLean B, Röst HL, Tate SA, Tsou C-C,
356 Reiter L, Distler U, Rosenberger G, Perez-Riverol Y, Nesvizhskii AI, Aebersold R, Tenzer S.
357 A multicenter study benchmarks software tools for label-free proteome quantification. Nature
358 Biotechnology 2016;34:1130.

359 **1 Material and methods**

360 **1.1 Sample preparation**

361 Lyophilized whole cell protein extracts of yeast and human were acquired from Promega and
362 lyophilized whole cell protein extract of ecoli was acquired from Waters. All extracts were already
363 reduced with dithiothreitol (DTT), alkylated with iodoacetamide and digested with Trypsin/Lys-C
364 Mix by their respective manufacturers. These extracts were reconstituted in 0.1% formic acid and
365 two master samples were created as in Navarro et al. [4], each in triplicate: A) a mixture of 65%
366 weight for weight (w/w) human, 15% w/w yeast and 20% w/w ecoli and B) a mixture of 65% w/w
367 human, 30% w/w yeast and 5% w/w ecoli. The resulting samples have logarithmic fold changes
368 (logFCs) of 0, 1 and -2 for respectively human, yeast and ecoli. One third of each of the six master
369 batches was mixed as a quality control (QC), resulting in ratios of 65% w/w human, 22.5% w/w
370 yeast and 12.5% w/w ecoli.

371 **1.2 Data acquisition**

372 For each of the six master samples three technical replicates injections were acquired to obtain nine
373 samples in total for both condition A and B. Nine technical replicate injections of the QC were also
374 acquired. All 27 samples were acquired in a randomized design in three different acquisition mass
375 spectrometry (MS) modes on three different mass spectrometers: 1) high definition MS^e (HDMS^e)
376 mode on a Synapt G2-Si (Waters), 2) data-dependent acquisition (DDA) mode on a Q-Exactive
377 (Thermo), and 3) in SWATH mode on a TripleTOF 5600 (AB Sciex) (Supplementary Figure 30).
378 For each of the $3 \cdot 9 \cdot 3 = 81$ samples, five μg was injected. All data was acquired in res mode.
379 The acquisition on the Synapt G2-Si was preceded by a nano-acquity (Waters) set up in microflow
380 liquid chromatography (LC), the acquistion on the Q-Exactive was preceded by a TODO micro LC,
381 an the acquistion on TripleTOF was preceded by an Eksigent micro LC. All samples were acquired
382 on a 150 minute gradient. After each three samples, an acoli autoQC sample was run to assess
383 the performance of the mass spectrometers. For the sequential window acquisition of all theoretical
384 mass spectra (SWATH) acquisition, TODO windows of TODO mass over charge ratio (m/z) were
385 used.

386 **1.3 Peak picking of raw data**

387 Raw data from all samples were peak-picked to obtain one comma separated values (csv) file per
388 sample in which all its ions, both low energy (LE) and high energy (HE), and intensities were
389 defined by their m/z apex, retention time (RT) apex, and drift time (DT) apex. In case of DDA or
390 SWATH, the DT apex is replaced by the m/z the precursor selection.

391 Waters' HDMS^e data was peak-picked with their Apex3D software, version 3.1.0.9.5 on a Windows
392 10 Workstation with 160 gigabytes TODO bits? (GB) random-access memory (RAM) and 16 central
393 processing units (CPUs). Selected parameters were a lockMass of 785.8426 for charge 2 with m/z
394 tolerance of 0.25, apexTrackSNRThreshold of 1, and write to Apex3D csv file instead of default
395 Apex2D csv file. Different counts thresholds of 1, 5, 10, 20, 50, and 100 were used for both LE and
396 HE to test the influence of noise on HistoPyA.

397 All resulting csv files were imported simultaneously in a Python environment to obtain a single
398 list containing all ions from all samples.

399 **1.4 Sample calibration and estimation of ion inter-run differences**

400 To calibrate the m/z , RT and optionally DT of each sample, all LE ions with an intensity larger than
401 2^{14} were selected and ordered by their m/z , regardless of sample origin. Between each consecutive
402 pair of ions, their m/z parts per million (ppm) error was calculated. Whenever a set of consecutive
403 ions, in which each sample was represented by exactly one ion, had smaller m/z ppm errors than
404 the left and right flanking m/z ppm errors, it was defined as a pseudo aggregate ion.

405 For each pseudo aggregate ion the point-to-point (ptp) distance in RT and optionally DT dimen-
406 sion of their representative ions was calculated. Based on the distribution of the median absolute
407 deviation of all RT or DT ptp errors, individual z -scores were calculated per pseudo aggregate ion.
408 Each pseudo aggregate ion with a z -score exceeding 5 was considered an outlier and removed. This
409 process of outlier removal was repeated until only pseudo aggregate ions with z -scores below 5 for
410 both their RT and DT remained.

411 50% of the pseudo aggregate ions were selected for calibration of the m/z and DT between each
412 sample. For each pseudo aggregate ion, the average RT, m/z , and DT was calculated. Per pseudo
413 aggregate ion the median ppm error of m/z and DT of all representative ions compared to the
414 pseudo aggregate ions average was calculated. These median sample ppm errors were subtracted
415 from the original m/z and DT of each ion present in the complete ion list. As a result, the median

416 error between all pseudo aggregate ions and the representative ions of each sample was zero.

417 The same 50% of pseudo aggregate ions were partitioned in groups to calibrate the RT between
418 samples. Two pseudo aggregate ions a and b belong to the same group if there exists a sample α in
419 which $RT_{a,\alpha} < RT_{b,\alpha}$ and a sample β in which $RT_{a,\beta} > RT_{b,\beta}$. Thus, two pseudo aggregate ions
420 c and d from two different groups always have representative ions so that for each sample γ the
421 statement $RT_{c,\gamma} < RT_{d,\gamma}$ is true. Per sample the average RT of each pseudo aggregate ion group
422 are taken as y -values, while the average RT of all representative ions of each pseudo aggregate ion
423 group are taken as x -values. Per sample, these x and y -values are then used to perform a piece-wise
424 linear transformation on the RT of all the ions in the complete ion list.

425 The remaining 50% of the pseudo aggregate ions were used to obtain an unbiased estimate of
426 the inter-run errors of the calibrated m/z , RT and DT errors. Per pseudo aggregate ion the ptp
427 distance (largest minus smallest of the representative ions) of the calibrated m/z , RT and DT
428 were calculated. The 99th percentile of each characteristic is now defined as the maximum allowed
429 inter-run error between two ions from different samples.

430 1.5 Ion inter-run alignment and noise definition

431 A network was created wherein each ion was a vertex. Between two ions an edge was set if and
432 only if the ions originated from different samples, were both acquired in either LE or HE, and had
433 calibrated m/z , RT and DT errors smaller than the maximum estimated inter-run errors.

434 Subsequently this network was trimmed, so that no path existed between two ions from the same
435 sample. This trimming was done iteratively on paths of increasing length. Whenever a path of
436 the specified length existed between two vertices from the same sample, all edges of the path were
437 removed. For each remaining connected components it was checked whether all ions originated from
438 different samples. If this was true, no further trimming happened on this connected component,
439 otherwise all edges which are not part of an edge-triangle are removed and the specified path length
440 was increased by one for the next trimming iteration.

441 The resulting network now consists of multiple connected components, in which each ion originates
442 from a different sample. Note that there may be connected components in which not all vertices are
443 connected, meaning that either some calibrated m/z , RT or DT exceed their respective maximum
444 allowed errors, or their connection got trimmed. The maximum allowed errors were determined on
445 the 99th percentile of pseudo aggregate ions, which in turn were defined with ions with intensity
446 above 2^{14} , meaning their apices were likely to be peak-picked more accurately than ions with lower

447 intensity. As such, these maximum allowed errors can be considered quite stringent and some
448 missing edges should be expected. Finally, each connected component was defined as an aggregate
449 ion. For all of these aggregate ions, their average calibrated m/z , calibrated RT and calibrated
450 DT was calculated. Each aggregate ion also has a weight that is defined by the number of samples
451 where it was detected. This property is proportional to the probability that this ion is a true signal.
452 Finally, all aggregate ions with only a single ions are considered noise and removed for subsequent
453 analyses.

454 To normalize intensity difference between samples, the average intensity of all aggregate ions
455 expressed in all samples was calculated, as well as the logFC distance of each individual sample to
456 this average. For each sample, the median of these logFC distances was determined and subsequently
457 subtracted from all ions in the complete ion list. Finally, the logFC of the average calibrated intensity
458 from ions in condition A compared to the average calibrated intensity from ions in condition B was
459 calculated per aggregate ion, or set to $-\infty$, null, $+\infty$ when no average could be calculated for
460 condition A and/or B.

461 **1.6 Estimation of intra-run differences between high energy aggregate ions of the 462 same precursor**

463 To estimate maximum RT and DT intra-run differences between aggregate ions derived from the
464 same precursor (e.g. fragments), HE isotopic aggregate ion pairs with ion representatives in all
465 samples are used. Two aggregate ions are defined as an isotopic pair if and only if their difference
466 in aggregate calibrated m/z is $1.002861 \pm x$ ppm (averagine isotope) with x the maximum inter-run
467 m/z error. Furthermore, the difference in original RT and DT per sample should be smaller than
468 the inter-run maximum error for each sample, assuming intra-run errors are smaller than inter-run
469 errors. Finally, this pair should be unique, meaning no other potential isotopic pair can be formed
470 with either of the aggregate ions. For this estimation, this generally implies only the mono-isotopic
471 and first isotope can be detected and that the second isotope is not present as an aggregate ion
472 expressed in all samples, or that a charge other than 1 was accidentally used.

473 Two ions from the same sample are now defined as co-eluting if and only if their distance in RT
474 and DT is smaller than the 99th percentile of the isotopic aggregate ion pair distribution per sample.

475 A special situation arises when determining co-elution between LE and HE scans for e.g. frag-
476 ments and precursors, as there is a drift shift between those channels. To correct this drift shift,
477 unfragmented pairs of fully reproducible LE and HE aggregate ions with equal m/z , within intra-

478 run ppm error, are determined in a similar way as isotopic pairs where original RT per sample
479 should be smaller than the inter-run maximum error for each sample. As with isotopic pairs, each
480 unfragmented pair should be unique. Hereafter, the relative drift shift, i.e. difference in drift time
481 divided by LE drift in ppm, per sample between LE and HE ions is determined and only those
482 within the 10th and 90th percentile are retained. Furthermore ions with DT below 50 or greater
483 than 190 are removed to avoid boundary issues. Optimal parameters a , b , c and d are then deter-
484 mined such that for all the retained ions the error between the relative drift shift y and the function
485 $y = a \cdot \|dt, mz\| + b \cdot \arctan(dt/mz) + c \cdot dt/mz + d$ has an optimal least squares fit. Finally, this
486 function is applied to all ions of all aggregate ions per sample.

487 TODO PPM difference calibration between HE-LE

488 **1.7 Aggregate ion network generation**

489 A network was created in which all aggregate ions were vertices. An edge is set between two
490 aggregate ions if and only if they consistently co-elute. Two aggregate ions are defined as *consistently*
491 *co-eluting* if and only if they co-elute for each overlapping sample. However, as the intra-run
492 differences within each sample are independent, a large sample count can introduce a dimensionality
493 curse, meaning it is unlikely that representative ions co-elute in each sample even if they originate
494 from the same precursor. Therefore the definition of *consistently co-eluting* is weakened to mean
495 that they should have a probability of at least 0.999 to overlap in at least x out of y samples.
496 Herein the probability is calculated by binomials, i.e. $\sum_{x \geq i}^y \binom{y}{i} \cdot 0.99^{y-i} \cdot 0.01^i > 0.999$. As a final
497 constraint, two aggregate ions should co-elute in at least two samples to be considered *consistently*
498 *co-eluting*.

499 **1.8 Fragment-centric annotation of an aggregate ion network**

500 At this point, the complete experiment has been collapsed into a single (noiseless) aggregate ion
501 network. Here, we used the aggregate ions for the final analysis, but this can easily be split into a
502 separate ion network for each sample. A fasta file containing all SwissProt entries from human, yeast
503 and ecoli was downloaded. The crap database was appended to this fasta, as well as a decoy with
504 all reversed protein sequences. A standard in silico tryptic digest with one miscleavage and default
505 amino acids masses, with the exception of a cysteine to which a carbamido mass of 57.021464 was
506 added, was made to obtain a list of peptides and their masses. Duplicate peptides from different
507 proteins were merged to obtain a list of unique peptide sequences. Peptides originating solely from

508 decoy proteins were classified as decoy peptides, while all others were classified as targets. For each
509 peptide, the masses of all b- and y-ions was calculated.

510 For each HE aggregate ion, all potential singly, doubly and triply charged b- or y-ion explanations
511 were determined within the inter-run ppm error. Moreover, each of these explanations belong to a
512 peptide, so every aggregate ion has a list of peptides from where it could have originated.

513 For each aggregate ion with at least three peptide explanations and at least two edges in the
514 aggregate ion network a hyperscore was determined in an X-Tandem! like fashion for all its potential
515 peptide explanations. For each of the peptide explanations it was counted how often it occurred in
516 the peptide explanations of the neighboring aggregate ions. Hereafter, the cumulative log frequency
517 of all but the highest of these counts was determined and used for a robust random sample consensus
518 (RANSAC) regression. A hyperscore equal to minus the regressed prediction of the highest count was
519 then determined for all peptides with this highest corresponding count. Note that some aggregate
520 ions have no peptides with a hyperscore, for instance when no regression can be made.

521 For each aggregate ion with at least one peptide with a hyperscore, it is checked whether there is
522 an LE aggregate ion that consistently co-elutes.

523 All aggregate ions and their remaining peptide explanations, meaning with a hyperscore and co-
524 eluting precursor, are now considered as a peptide-ion match (PIM) and given to percolator where
525 they are treated as if they were peptide-spectrum matchs (PSMs). Percolator features are set to RT,
526 fragment delta mass (ppm), precursor delta mass (ppm), neighbor count, peptide count, hyperscore,
527 precursor charge and fragment ion type with e.g. b7 as 7 and y4 as -4. Percolator was run with
528 default parameters with the addition of post processing tdc (Y-flag) to correct for an imbalance in
529 targets and decoys, and all predicted features (D-flag set to 15). Finally, all PIMs with a *q*-value
530 below 0.01 are retained.

531 For each of the aggregate ions belonging to a PIM with *q*-value below 0.01, an exhaustive anno-
532 tation is done for all its neighbors, which can thus be annotated as singly, doubly or triply charged
533 precursor, b-NH3, b-H2O, c, a, a-NH3, a-H2O, y, y-NH3, y-H2O or x fragment of a specific peptide.

534 **1.9 AutoQC, data deposition and source code availability**

535 TODO

536 **References**

- 537 [1] Limonier F, Willems S, Waeterloos G, Sneyers M, Dhaenens M, Deforce D. Estimating the
538 reliability of low-abundant signals and limited replicate measurements through MS2 peak area
539 in SWATH. PROTEOMICS;0.ja:1800186.
- 540 [2] Craig R, Beavis RC. TANDEM: matching proteins with tandem mass spectra. Bioinformatics
541 2004;20.9:1466–1467. eprint: /oup/backfile/content_public/journal/bioinformatics/
542 20/9/10.1093/bioinformatics/bth092/2/bth092.pdf.
- 543 [3] The M, MacCoss MJ, Noble WS, Käll L. Fast and Accurate Protein False Discovery Rates on
544 Large-Scale Proteomics Data Sets with Percolator 3.0. Journal of The American Society for
545 Mass Spectrometry 2016;27.11:1719–1727.
- 546 [4] Navarro P, Kuharev J, Gillet LC, Bernhardt OM, MacLean B, Röst HL, Tate SA, Tsou C-C,
547 Reiter L, Distler U, Rosenberger G, Perez-Riverol Y, Nesvizhskii AI, Aebersold R, Tenzer S.
548 A multicenter study benchmarks software tools for label-free proteome quantification. Nature
549 Biotechnology 2016;34:1130.

550 **Supplementary**

551 **Peak picking**

Table 1: Apex 3D peakpicking results on the 27 LFQ samples.

552 Schematic overview of quick alignment to detect pseudo aggregate ions

Pseudo aggregate formation

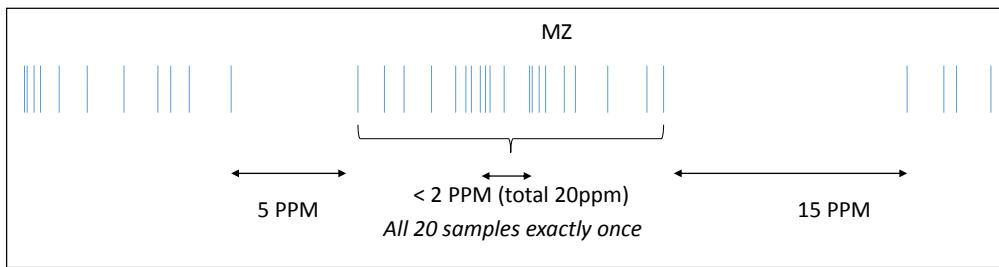


Figure 7: Schematic overview of quick alignment to detect pseudo aggregate ions.

553 Schematic overview of pseudo aggregate ion group piece-wise linear transformations
554 for retention times

RT Calibration

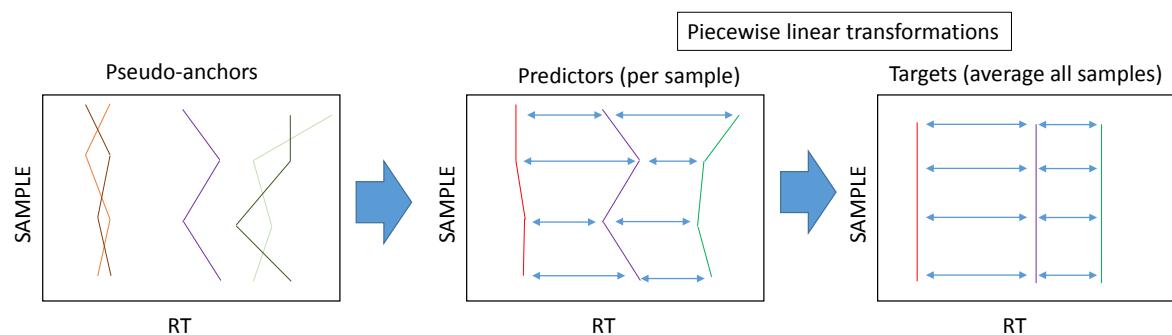


Figure 8: Schematic overview of pseudo aggregate ion group piece-wise linear transformations

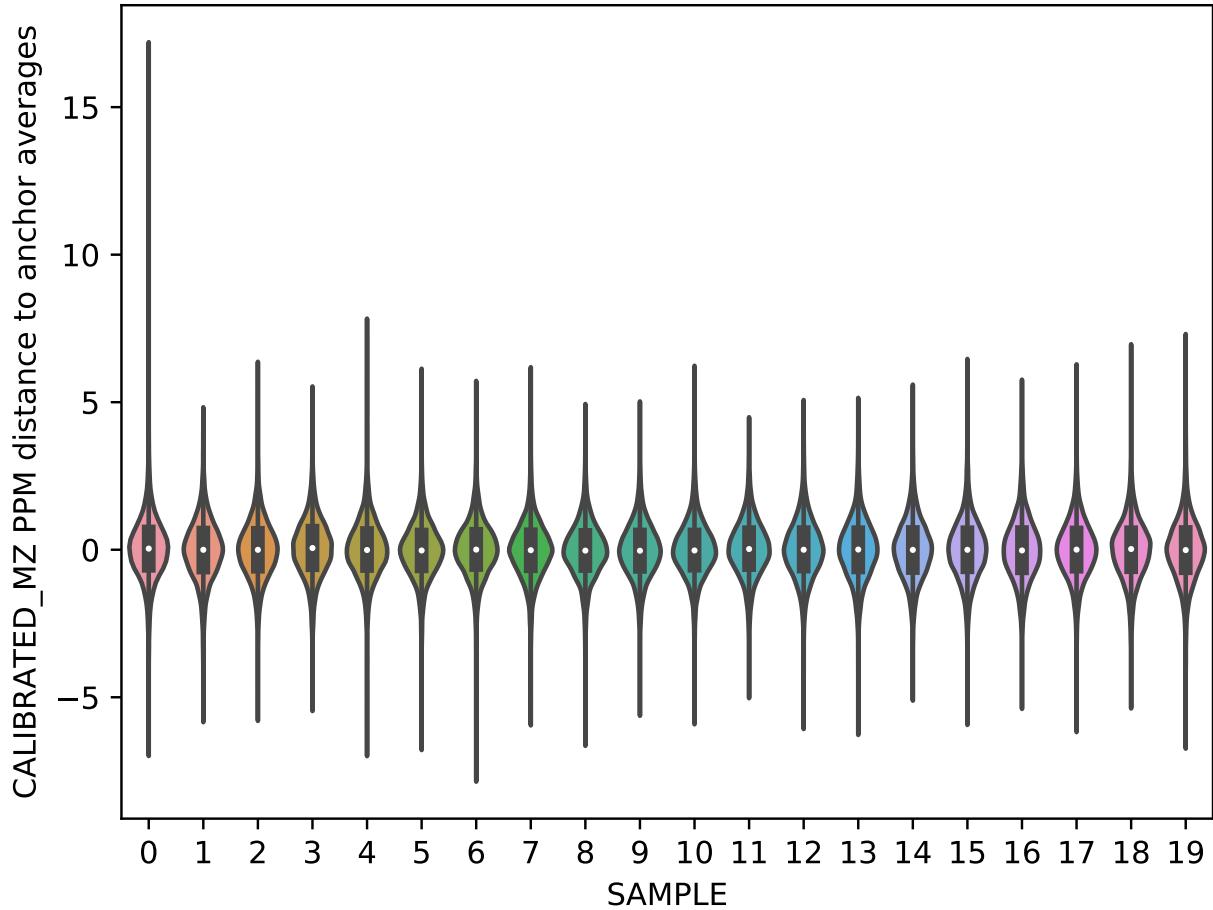


Figure 9: Estimation of m/z error per sample toward the pseudo aggregate ions.

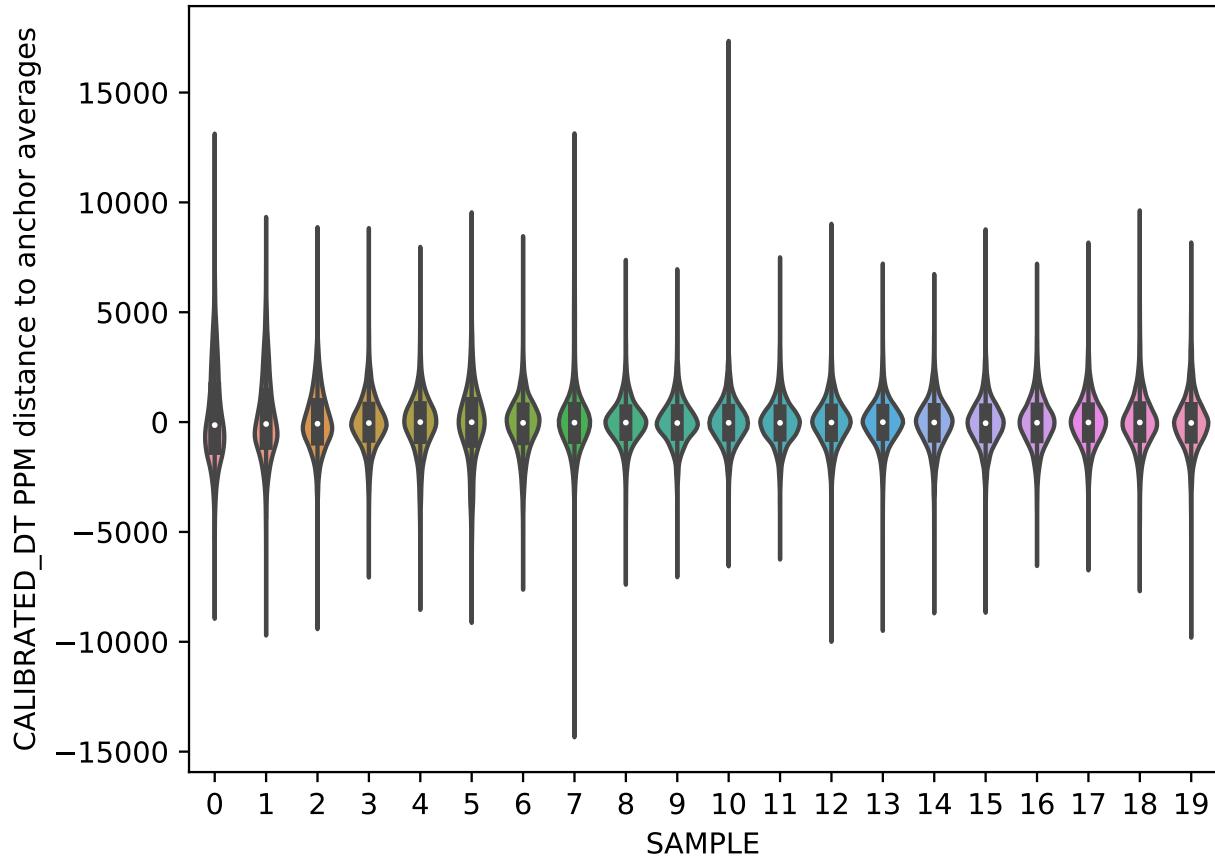


Figure 10: Estimation of DT error per sample toward the pseudo aggregate ions.

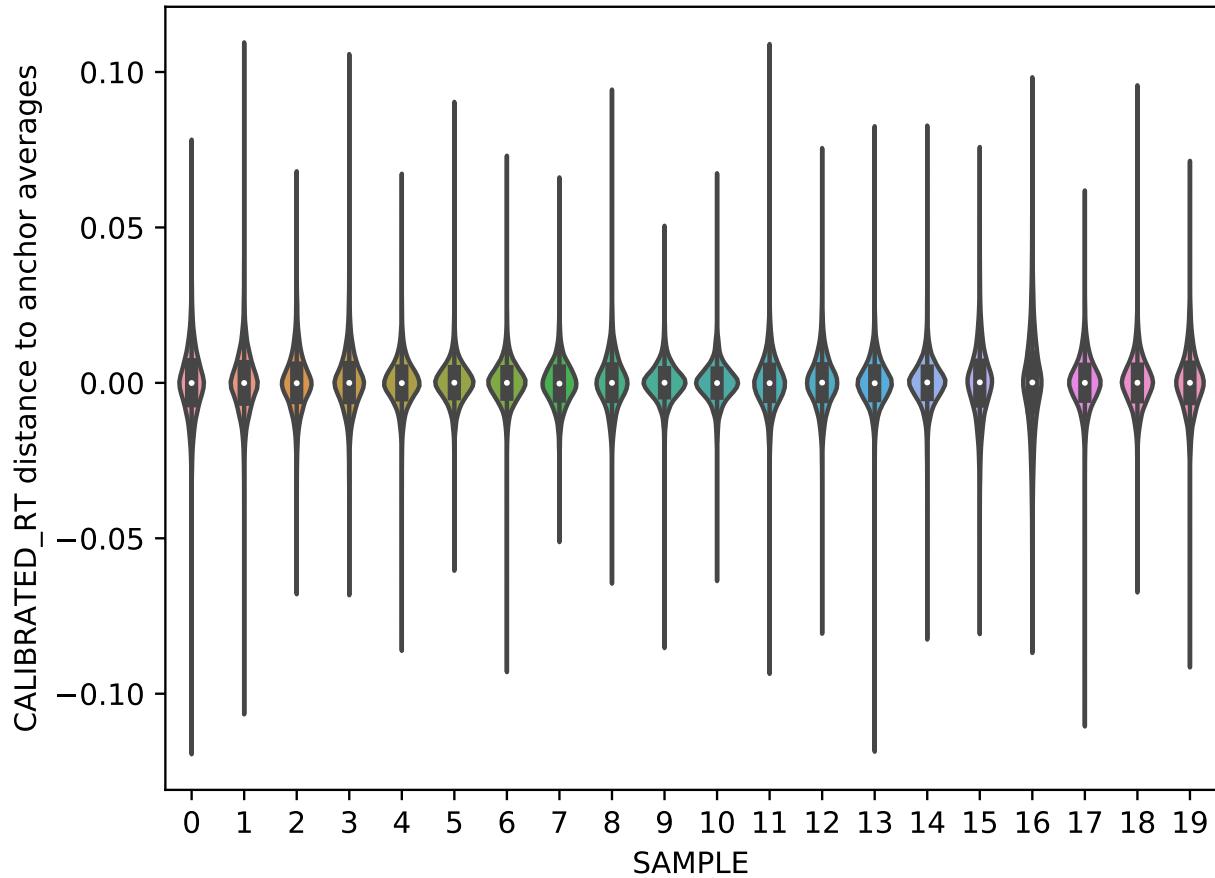


Figure 11: Estimation of RT error per sample toward the pseudo aggregate ions.

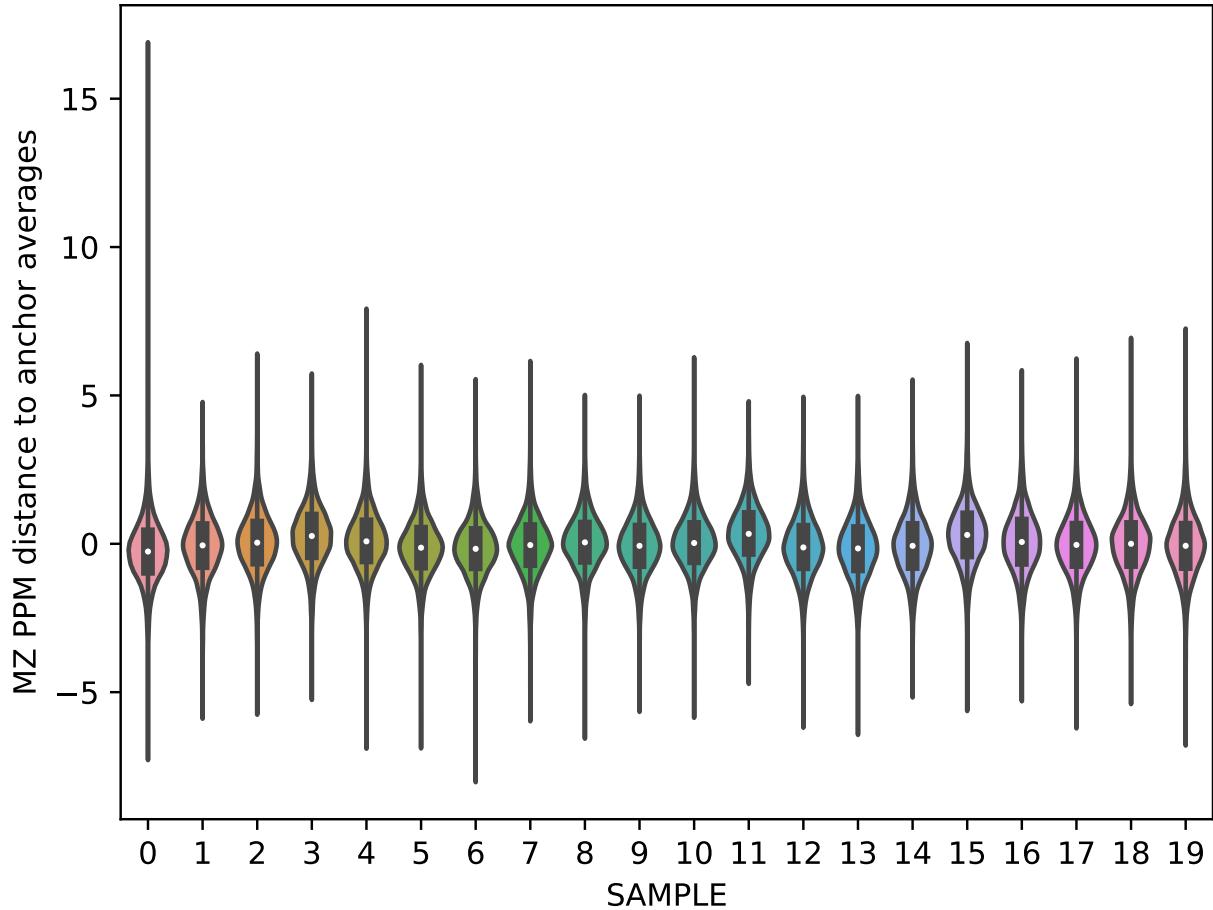


Figure 12: Uncalibrated m/z error per sample toward the pseudo aggregate ions.

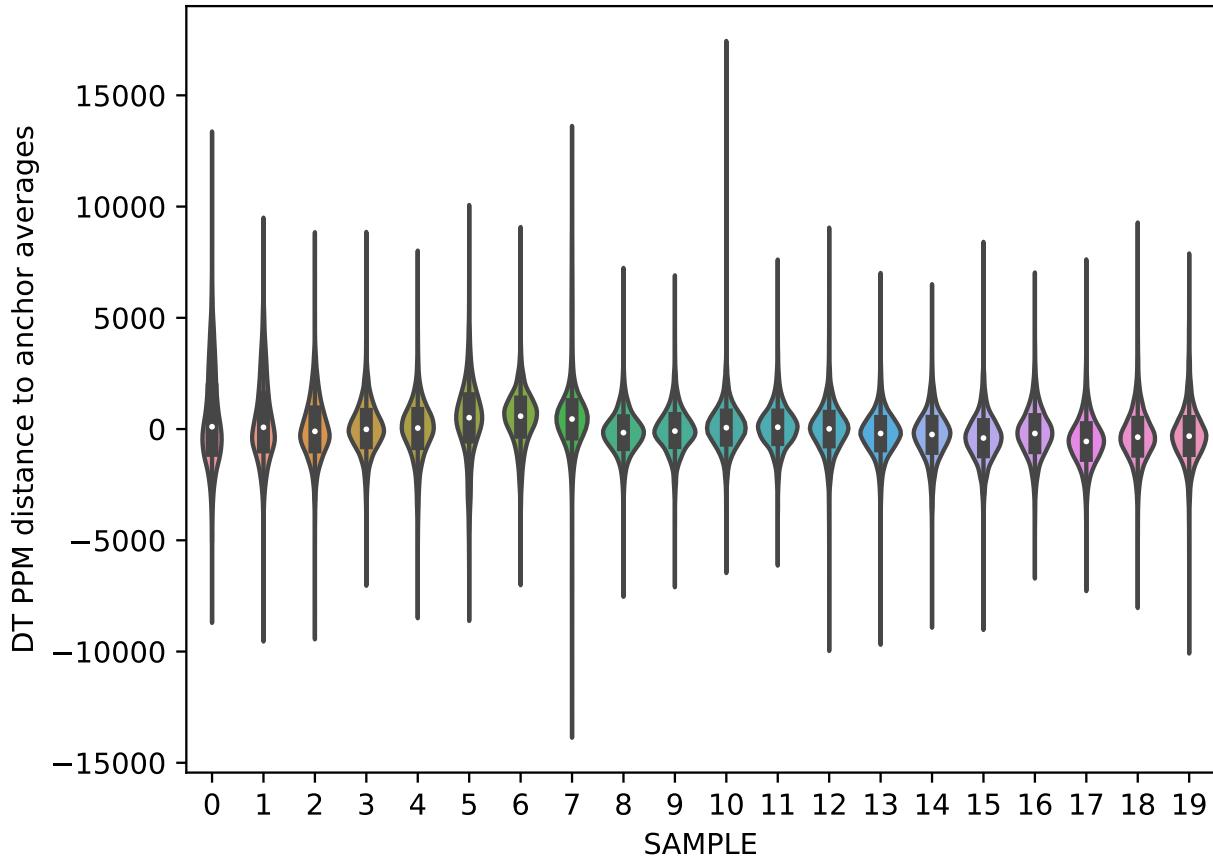


Figure 13: Uncalibrated DT error per sample toward the pseudo aggregate ions.

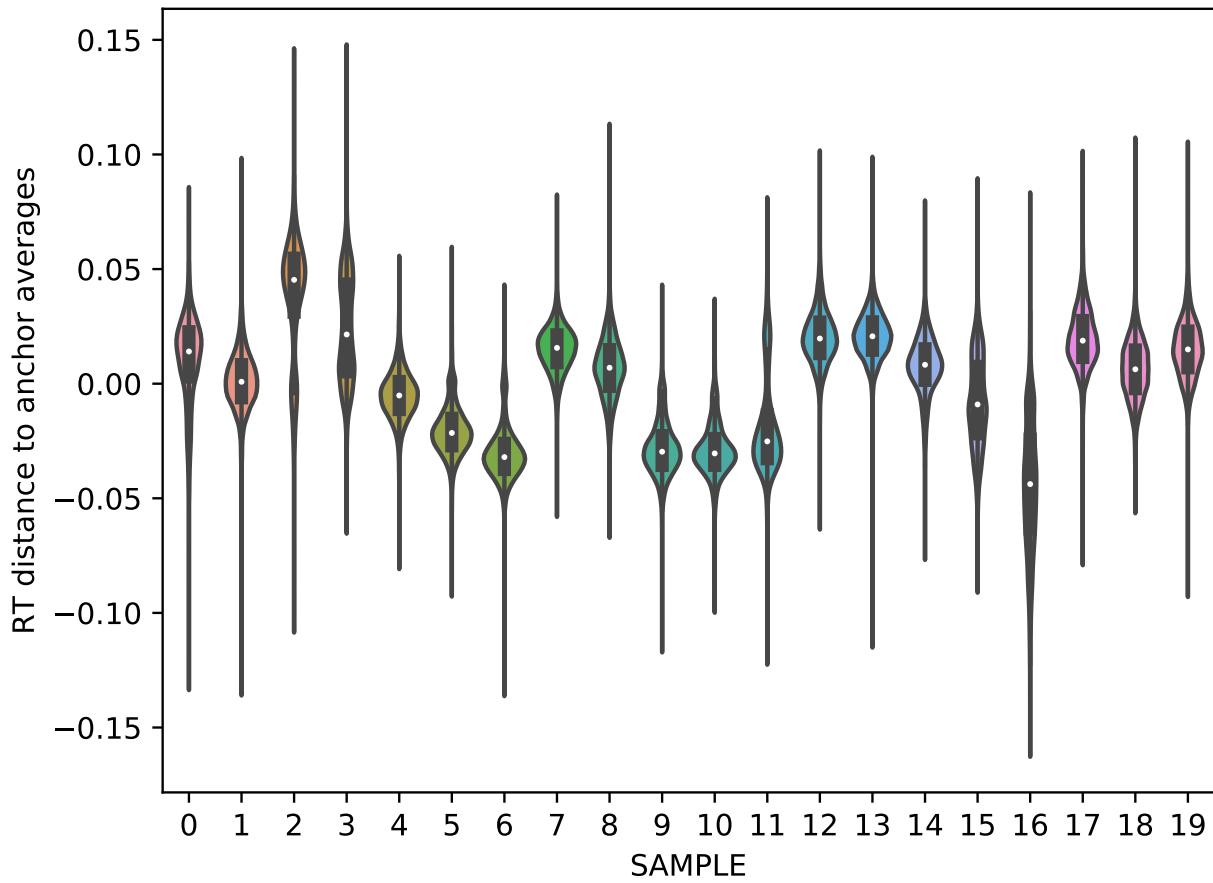


Figure 14: Uncalibrated RT error per sample toward the pseudo aggregate ions.

561 **Low energy aggregate ion sizes**

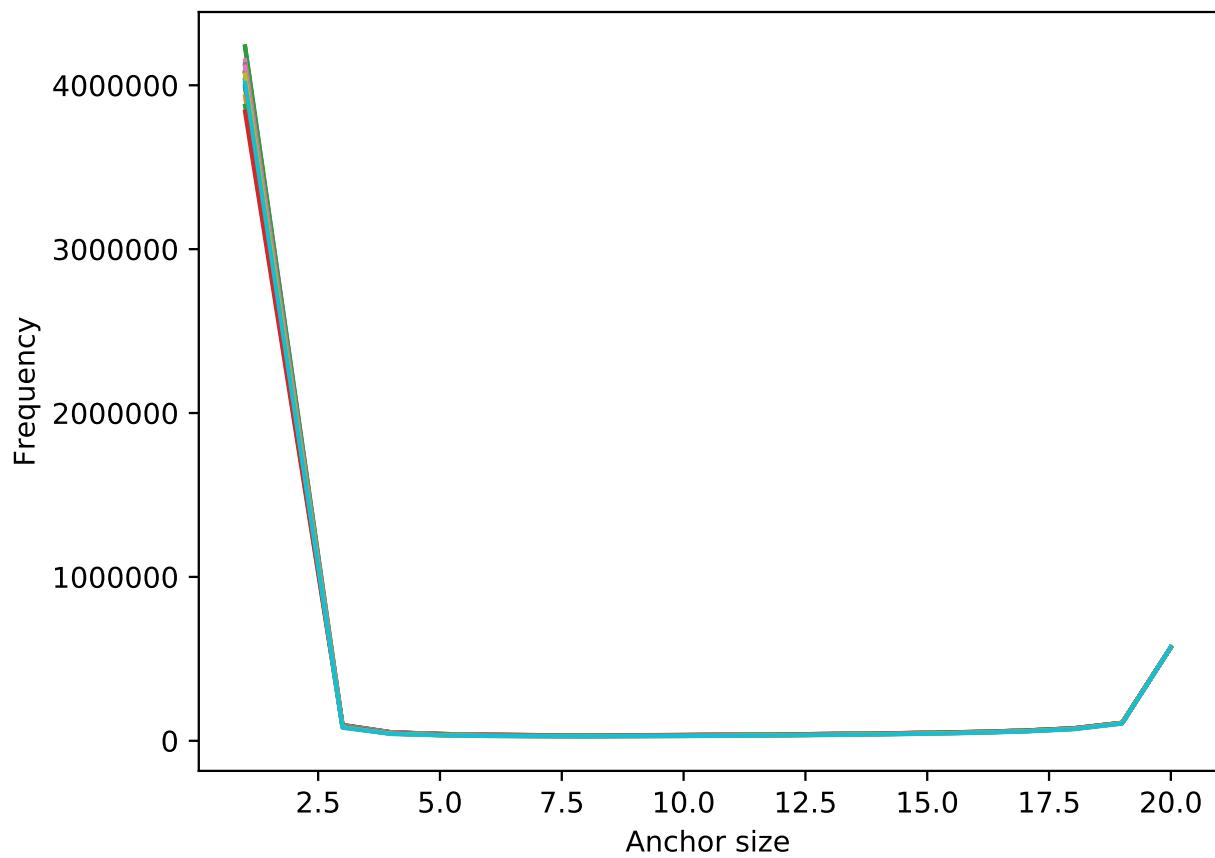


Figure 15: Low energy aggregate ion sizes.

562 **High energy aggregate ion sizes**

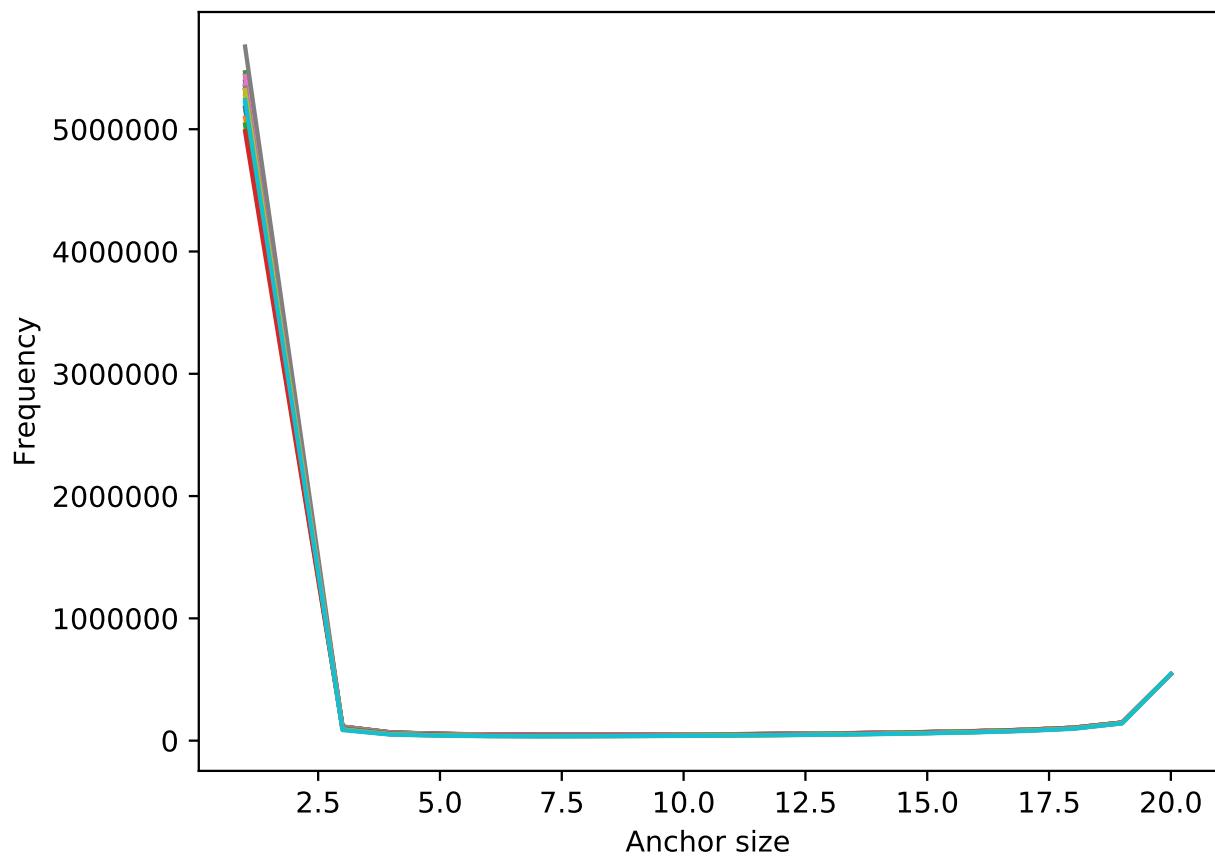


Figure 16: High energy aggregate ion sizes.

563 Low energy aggregate ion intensities

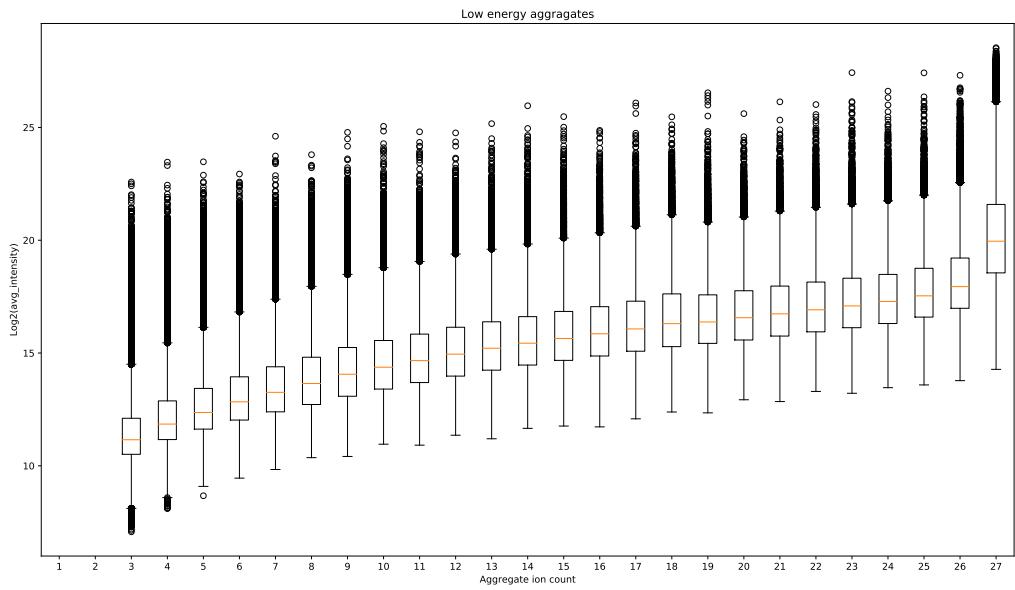


Figure 17: Low energy aggregate ion intensities.

564 **High energy aggregate ion intensities**

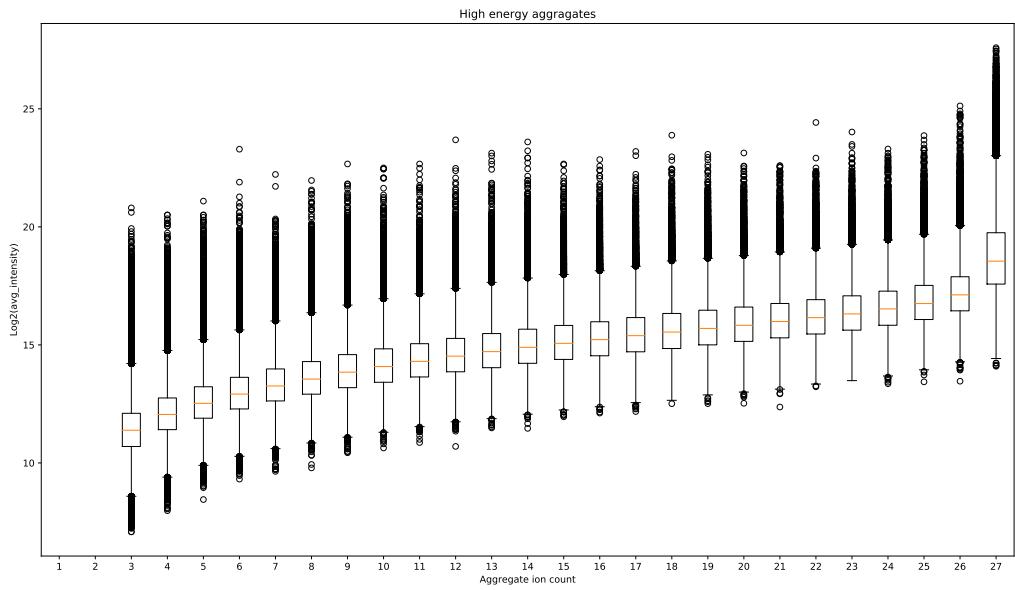


Figure 18: High energy aggregate ion intensities.

565 CV of uncalibrated intensities in condition A

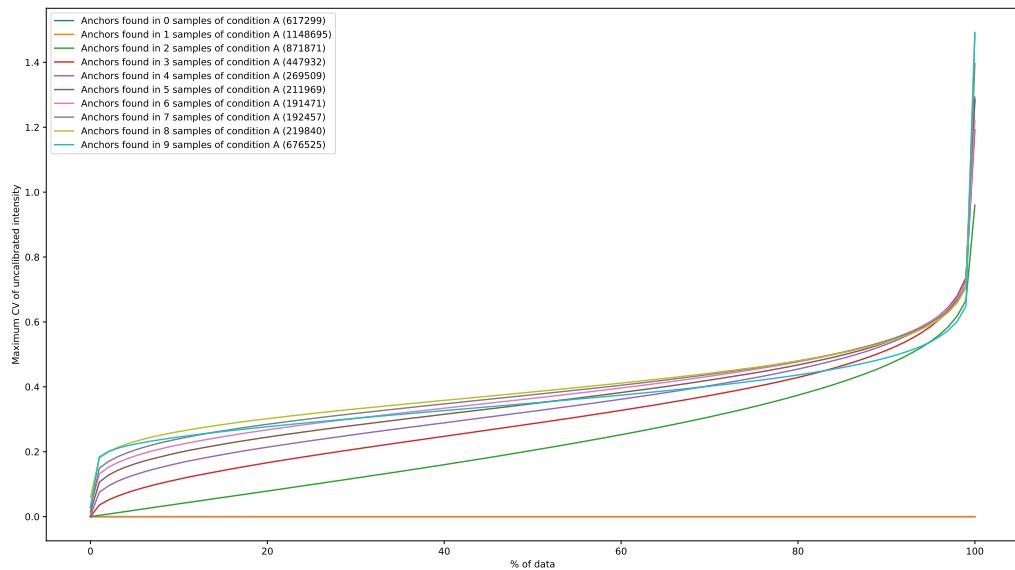


Figure 19: CV of uncalibrated intensities in condition A.

566 CV of uncalibrated intensities in condition B

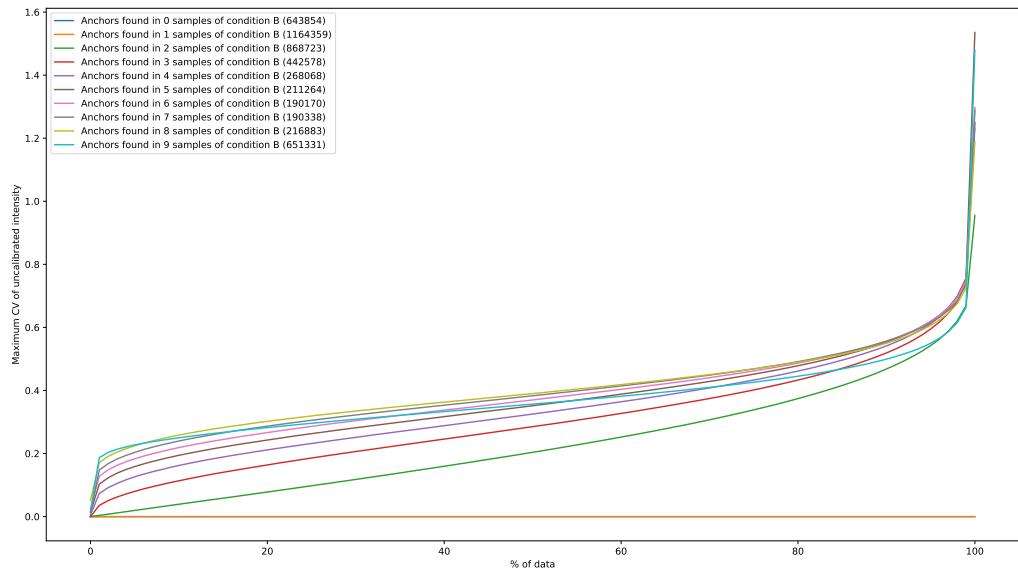


Figure 20: CV of uncalibrated intensities in condition B.

567 CV of uncalibrated intensities in condition QC

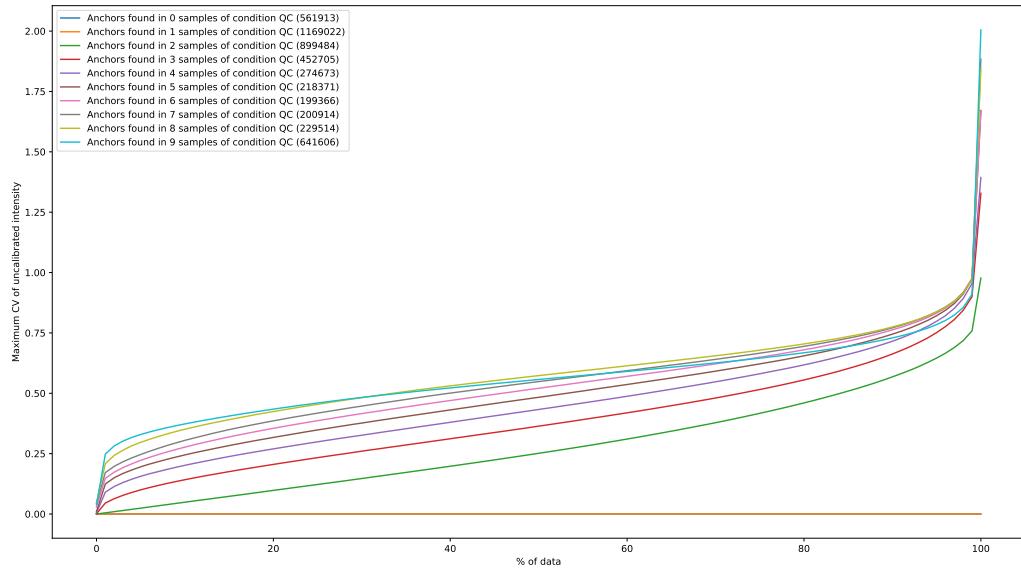


Figure 21: CV of uncalibrated intensities in condition QC.

568 CV of calibrated intensities in condition A

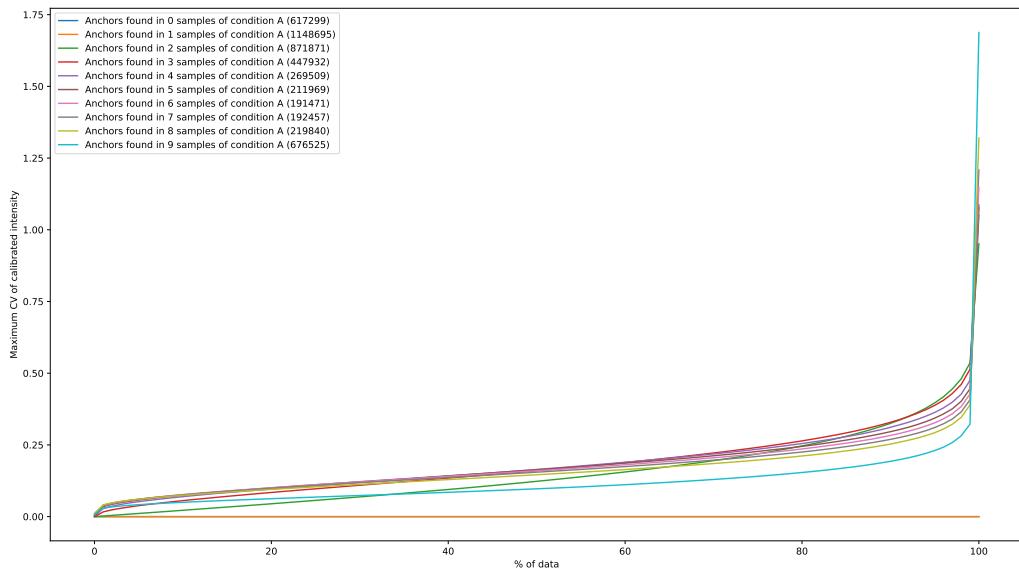


Figure 22: CV of calibrated intensities in condition A.

569 CV of calibrated intensities in condition B

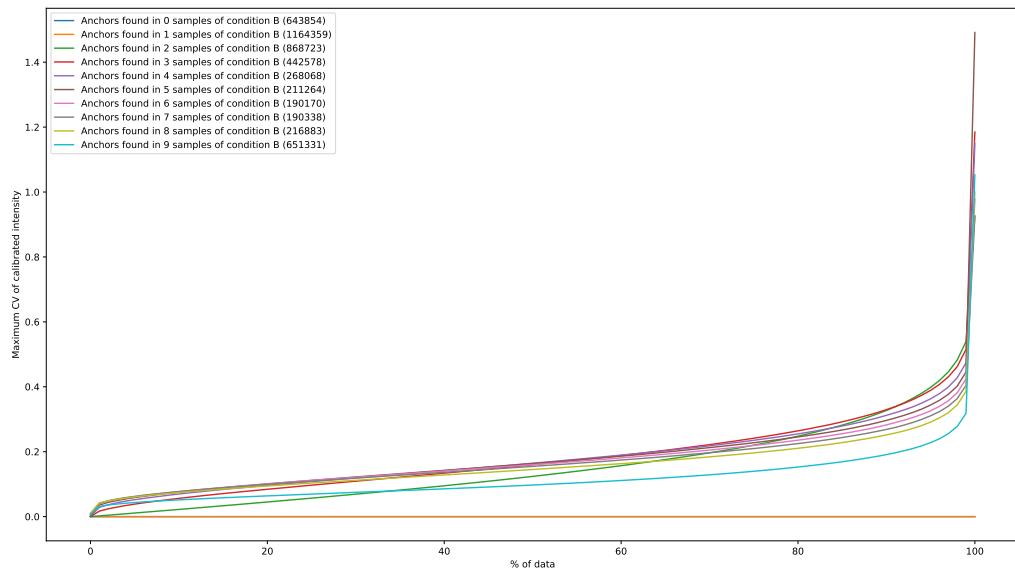


Figure 23: CV of calibrated intensities in condition B.

570 CV of calibrated intensities in condition QC

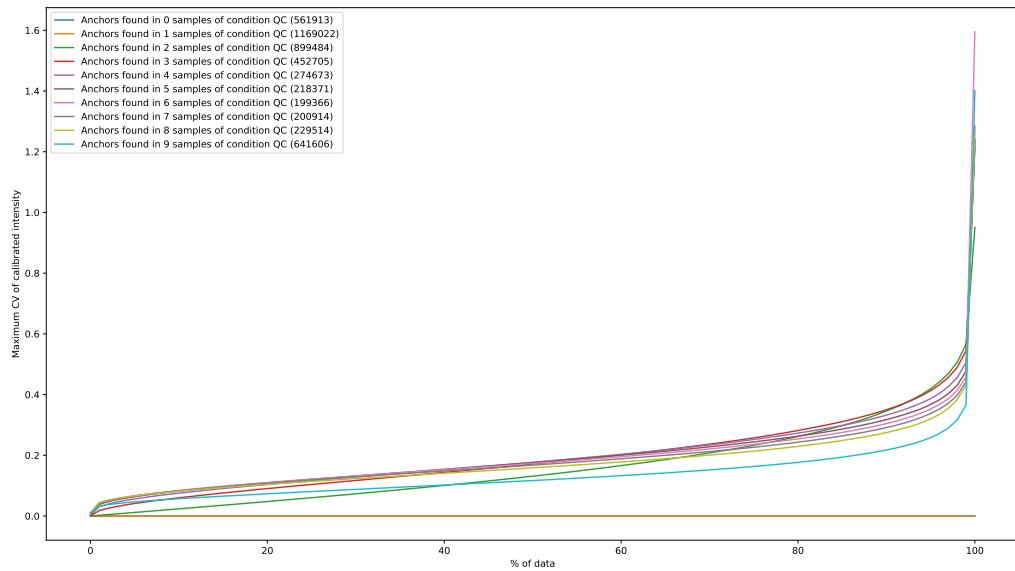


Figure 24: CV of calibrated intensities in condition QC.

571 **Distribution of logarithmic fold changes**

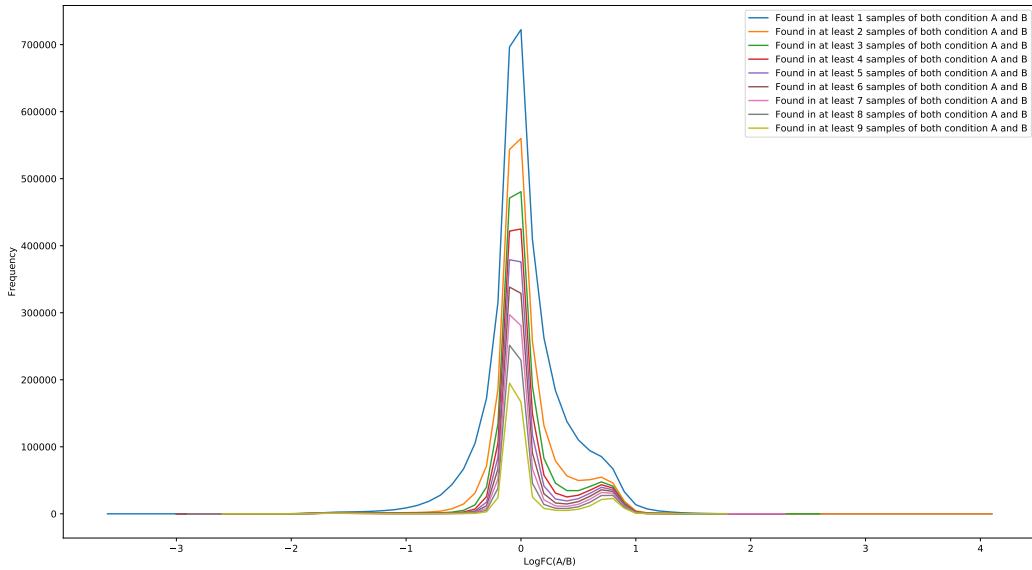


Figure 25: Distribution of logarithmic fold changes.

Signal detection

	Cumulative	Per anchor
>>> for i in zip(a1, b1, b5, b10, b20): ("{:>15})*len(i)).format(*i)		
... 1 192802258 128543866 95512773 55882459'		
3 5794722 5771011 5648147 5137175'		
4 4512293 4594136 4445254 4161202'		
5 397711 397711 395244 3720718'		
6 3627869 3624534 3594377 3428395'		
7 3369474 3369924 3342721 3283782'		
8 3157896 3155981 3136182 3010576'		
9 2953119 2951119 2926878 2820498'		
10 2807782 2806540 2792718 2760778'		
11 2651562 2650510 2638849 2557968'		
12 2507260 2506102 2494102 2424844'		
13 2355116 2354406 2346207 2284544'		
14 2209839 2209203 2205254 2140312'		
15 2062699 2062215 2056736 2011825'		
16 1925359 1925359 1906652 1867460'		
17 1752642 1752495 1749275 170552'		
18 1586616 1586495 1578354 1557323'		
19 1383131 1383055 1381798 136818'		
20 1115224 1115192 1114630 1106263'		
>>> for i in zip(a1, b1, b5, b10, b20): ("{:>15})*len(i)-np.diff(i[1:])).format(*i)		
... 64324681 3299729 39197842'		
1599 63951 2369509'		
1578 13734 14926'		
1519 10816 43144'		
765 5974 27643'		
611 4484 18113'		
375 3274 15294'		
322 2669 12382'		
198 2161 11869'		
203 1085 9933'		
219 1498 9315'		
84 1510 8451'		
152 1256 8301'		
69 1356 8095'		
97 1093 8893'		
98 1015 7692'		
71 859 7513'		
44 695 7113'		
32 562 6367'		
>>> for i in zip(a1, b1, b5, b10, b20): ("{:>15})*len(i)-np.diff(i[1:])).format(*i)		
... 64324681 3299729 39197842'		
4079 15711 68959'		
13112 74936 399794'		
7595 50888 215728'		
4599 35844 162258'		
4277 11968 159318'		
3088 20192 121632'		
2898 24821 111438'		
1969 21010 109309'		
22333 20724 108933'		
2628 17976 111780'		
1092 19639 169963'		
12128 158569 155114'		
968 18840 121425'		
1552 17488 129488'		
1539 17255 138764'		
1278 15444 131510'		
836 13205 135347'		
640 11240 127340'		

- Differences (anchor count)

- Differences (ion count)

- Ion Alignment

- Different peak picking thresholds
- K562 anchors
 - Peak picking 1,5, 10, 20
 - ?Robust to noise
 - ?CVs validation

Figure 26: HistoPyA’s robustness to peak picking noise.

573 Co-elution retention time distances

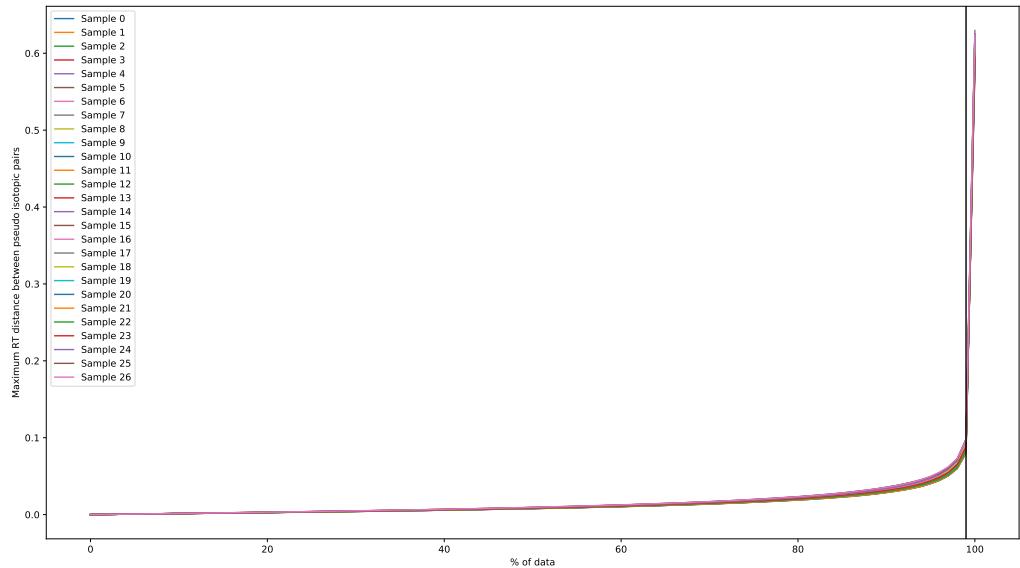


Figure 27: Distribution of intra-run retention time distances.

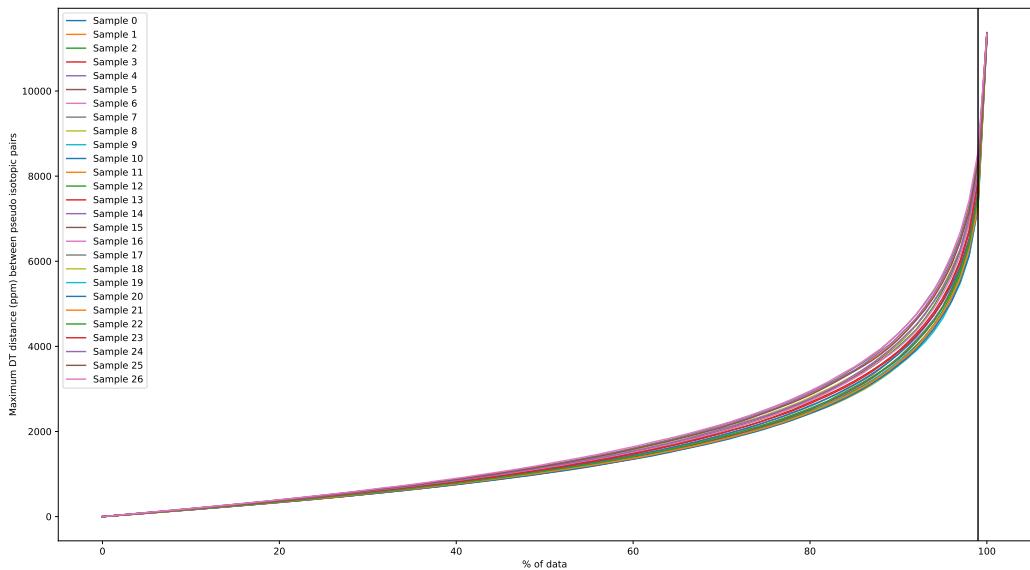


Figure 28: Distribution of intra-run drift time distances.

575 **Neighbors distribution**

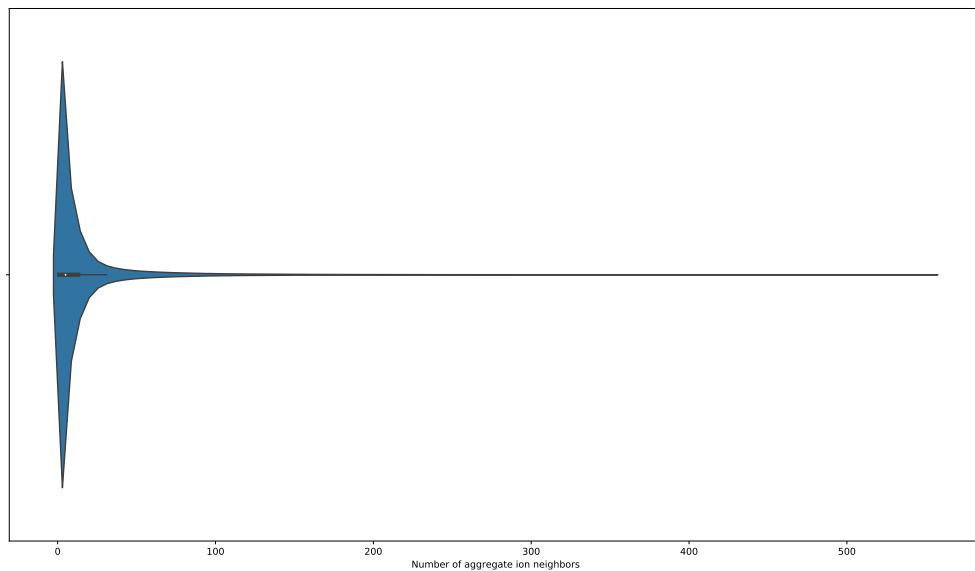


Figure 29: Distribution of aggregate ion neighbor counts.

576 **Sample scheme**



Figure 30: Sample scheme.

577 **Drift time shift**



Figure 31: Drift time shift between low energy and high energy channels.