

Single window ion mobility data independent acquisition mass spectrometry of multiple samples allows fragment-centric annotation of near-noiseless ion-networks

Sander Willems¹, Simon Daled¹, Bart van Puyvelde¹, Laura de Clerck¹, Filip van Nieuwerburgh¹, Dieter Deforce¹, and Maarten Dhaenens¹

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

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

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

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

26 1 Introduction

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

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

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

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

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

Here, we present single window ion mobility (SWIM)-DIA. Herein a sample is put on an

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

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.

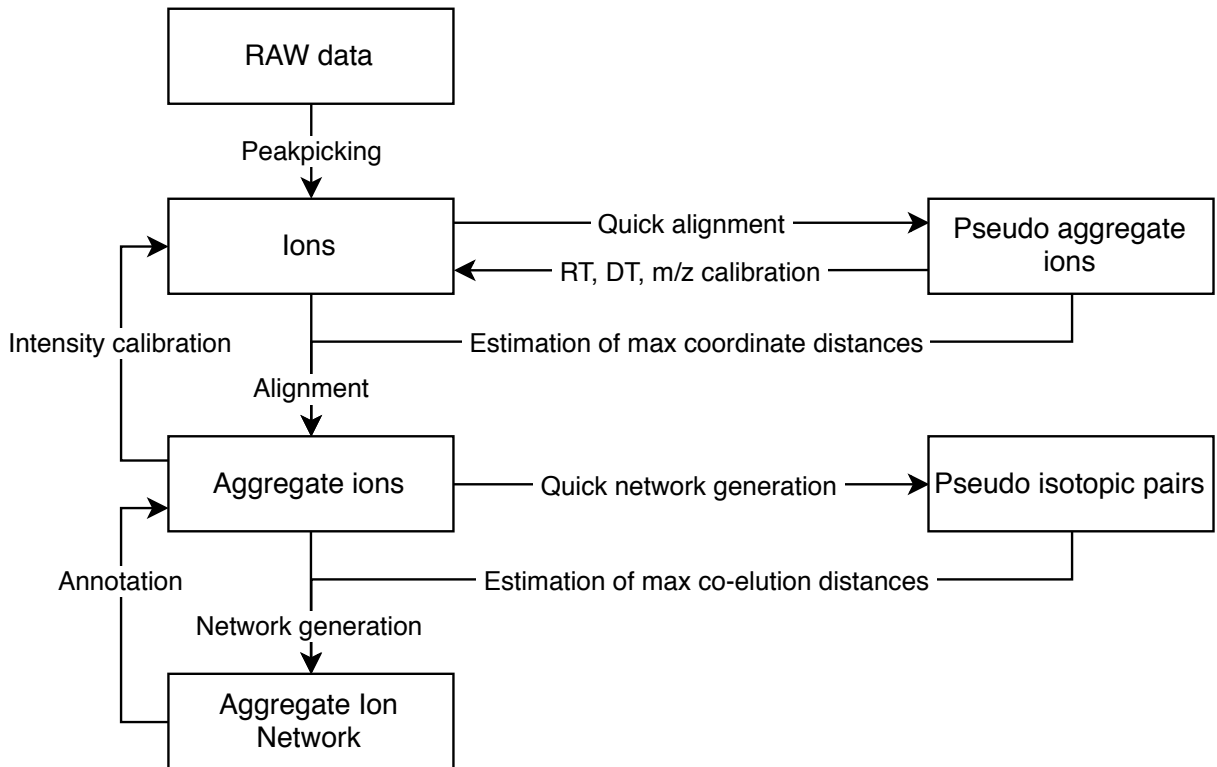


Figure 1: Schematic overview of HistoPyA's workflow.

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 2 Results

128 2.1 Peak picking and inter-run ion alignment

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

132 Using Water’s Apex 3D peak picking algorithm the raw data from 27 samples (nine per condition
133 A, B and QC) were peak picked in all of its dimensions: m/z , RT, DT and intensity. Even though
134 these samples had different ratios for human, yeast and ecoli proteins, the resulting number of
135 detectable ions was similar for all of them, with each sample resulting in roughly 6 million ions in
136 both LE and HE ions (Supplementary Table 1) at the lowest thresholds of 1 for both LE and HE.

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

To compensate for small differences in acquisition, all 27 samples were calibrated with a quick alignment. To perform this quick alignment the peak picked ions of all samples were merged into a single list and only those ions with an intensity larger than 2^{14} were used. Of all these ions more than 10 thousand pseudo aggregate ion could be formed as they had a unique m/z which was found exactly once in each sample (Supplementary Figure 4). A thousand of these pseudo aggregate ions had either RT or DT outliers and were removed, and the remaining pseudo aggregate ions were equally partitioned in a group for calibration and a group for estimation. With the pseudo aggregate ions for calibration both the m/z and DT of all peak-picked ions were corrected by at most 10 parts per million (ppm) and 10000 ppm respectively, depending on their sample origins. For the RT calibration the pseudo aggregate ions were grouped in more than 500 groups so that each group had a distinct RT and the RT of each ion of each pseudo aggregate ion in the group was strictly smaller than the RT of each ion in the next pseudo aggregate ion group. With these pseudo aggregate ions groups a piece-wise linear transformation was performed to calibrate the ions of each sample, resulting in a more consistent retention time (Supplementary Figure 5).

After this calibration, the other half of the pseudo aggregate ions were used to estimate the maximum inter-run distance for each of RT, m/z , DT, resulting in respectively 5 ppm, 6000 ppm and 0.2 minutes (Supplementary Figures 6, 7 and 8). These estimates are relatively small, indicating that the calibration performed well or that the acquisition was very robust. While the latter is subjective to evaluate without comparison, the estimates with uncalibrated coordinates are more than twice as large (Supplementary Figures 9, 10 and 11), indicating that at least the former statement is true.

2.1.2 Inter-run ion alignment

With the estimation parameters obtained from the pseudo aggregate ions, more than 100 million pairs of ion neighbors could be defined. However, the definition of neighboring implies a transitive relation over a path of multiple neighboring pairs. With this transitivity, there are clusters with over a 100 ions in them, sometimes with a sample being represented 10 times. To avoid such ambiguous situations, the pairs of ion neighbors were trimmed so that each remaining cluster, hereafter defined as an aggregate ion, has at most one ion per sample. With this trimming 10 million pairs of ion neighbors were removed. This relatively low percentage of trimming implies the estimation parameters are stringent enough to obtain a good specificity. Most of the resulting aggregate ions

169 were either fully reproducible or random noise (Supplementary Figures 12 and 13). As expected,
170 the fully reproducible aggregate ions have the highest average intensities, while the noisy aggregate
171 ions are less intense (Supplementary Figures 14 and 14).

172 2.1.3 Intensity calibration and validation

173 Intensity was not used to create the aggregate ions and was therefore not calibrated with the pseudo
174 aggregate ions. This calibration is done after the full alignment to include as much information as
175 possible. On average, this calibration reduced the intensity CV of anchors by a factor 3 (Supple-
176 mentary Figures 16, 17, 18, 19, 20 and 21), resulting in e.g. more than 80% of all fully reproducible
177 anchors with a CV below 20 for each of condition A, B and QC.

178 With this calibrated intensity, the logFC values between condition A and B could be determined.
179 As expected, these logFC values indeed show three groups of -2, 0 and 1, corresponding to the
180 original mix ratios of the samples (Supplementary Figure 22).

181 2.1.4 Robustness to noise

182 To assess HistoPyA's ability to deal with noise, different peak picking thresholds in Waters' Apex3D
183 software were tested. For nearly all partially reproducible aggregate ions, additional samples could
184 be found at lower ion count thresholds. Fully reproducible aggregate ions were not hindered by the
185 potential extra interference (Supplementary figure 23). However, halving the ion count threshold in
186 Apex3D results in doubling the amount of noisy ions in HistoPyA. Nearly all ions are reproducible
187 signal at an ion count threshold of 100.

188 2.2 Aggregate ion network

189 With the 5 million aggregate ions obtained at a peak picking ion count threshold of 1, a quick
190 assessment of potential isotopic pairs was made. From the 1 million fully reproducible aggregate
191 HE ions, 100 thousand pairs of pseudo isotopic pairs were obtained. Based on these isotopes, the
192 maximum difference in RT and DT peaks was estimated to be smaller than 0.05 minutes and 7000
193 ppm respectively for all of the samples (Supplementary Figures 24 and 25).

194 With these maximum difference per sample the potential *consistent* co-elution between all HE
195 aggregate ions was determined. Hereby, a total 400 million pairs of aggregate ions could be defined
196 as *consistently* co-eluting, implying a probable same precursor in LE. While there are some aggregate

ions with more than a 100 neighbors, most aggregate ions have between 1 and 13 neighbors showing a high specificity (Supplementary Figure 26). Of note, 15% of all anchor ions do not have a single neighbor and could be considered noise, even though they are reproducible over multiple runs.

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

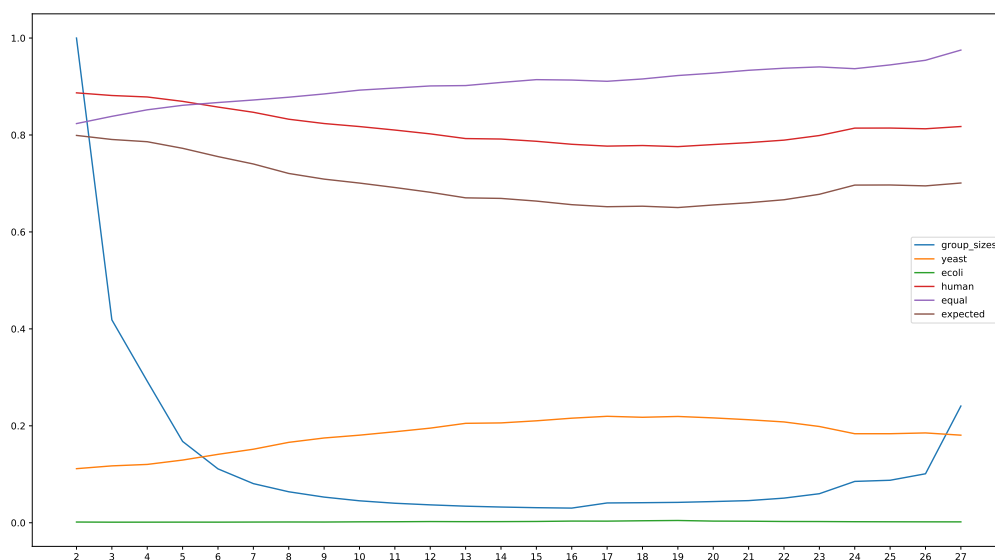


Figure 2: Neighbor organism origins.

When HE aggregate ions are deisotoped based on their neighbors, the delta mass frequencies are most prevalent for amino acid delta masses (Supplementary methods TODO). When only LE aggregate ions are considered and deisotoped, the singly charged aggregate ions have drift time that confirm the charge state as being 1.

To determine consistent co-elution between LE and HE, only consistency in RT was considered. This was essential because there is a DT shift between LE and HE scans of roughly 3 (Supplementary Figure 28)

Poor intensity correlation between neighbors (CV dependant)?

215 **2.3 Annotation**

216 **2.4 Model validation**

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

218 Fragment-centric annotation of DDA data gives similar results as Mascot annotation.

219 **2.5 Results**

220 Many fragments were significant and unique.

221 Overview of fragment features (used for percolator).

222 Percolator increased significance by many percent.

223 Many ions in aggregate ion network were annotated by significant proxies, ie limited "useless"
224 ions, ie noiseless.

225 Precursor existence filter very useful, but not essential.

226 More DIA peptides found than DDA.

227 Vendor specs of >12000 peptides for human reached.

228 **2.6 Annotation validation**

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

231 Venn diagrams of mascot and histopya peptides / proteins.

232 MRM validation of novel peptides not found in DDA?

233 RTs coincide with mascot DDA rts.

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

235 DTs coincide for same peptides: charge states and in source fragmentation confirmation.

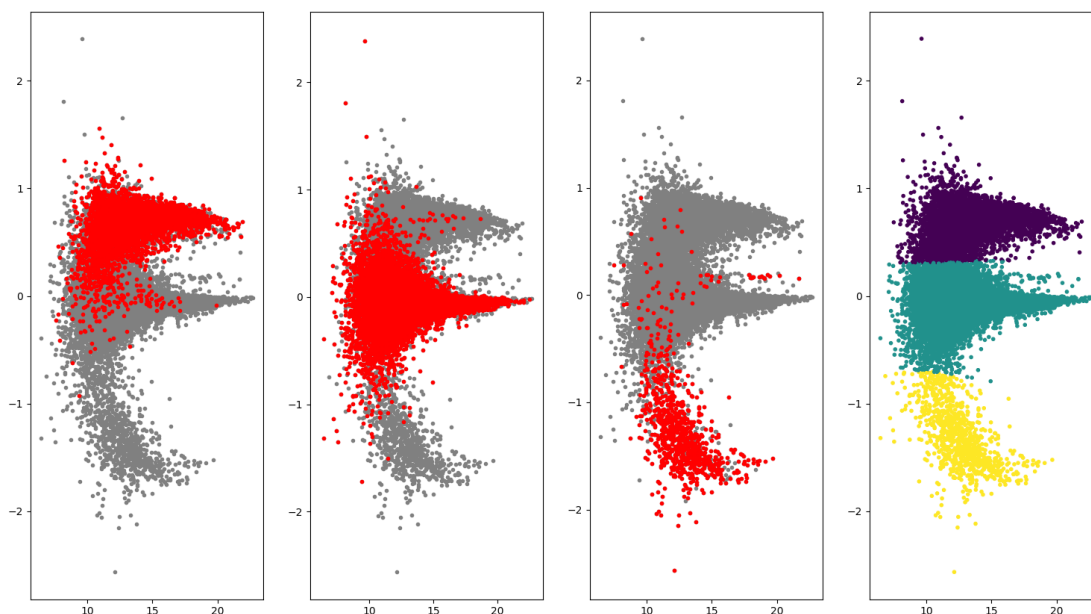


Figure 3: Experimental annotation and theoretical organism classification

3 Discussion

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

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

Cosmology vs quantum mechanics

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

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

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

HE only could allow better/simpler/cheaper instrumentation.

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

Data and software availability

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

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

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

Conflict of interest

The authors declare no conflict of interest.

Acknowledgements

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

The authors would like to express their gratitude to the Waters informatics team (including, but not limited to; Hans Vissers, Scott Geromanos and Steve Cievarini) and Lennart Martens (Ghent University (UGhent)) for their critical feedback throughout the project. Samples were acquired at the ProGenTomics facility and technical mass spectrometry assistance was provided by Sofie vande Castele (UGhent). Computational assistance was provided by Yannick Gansemans (UGhent) and Laurentijn Tilleman (UGhent).

272 **Author contributions**

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

278 **Acronyms**

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

280 **CCS** collisional cross section

281 **CPU** central processing unit

282 **csv** comma separated values

283 **CV** coefficient of variation

284 **DDA** data-dependent acquisition

285 **DIA** data-independent acquisition

286 **DT** drift time

287 **DTT** dithiothreitol

288 **ESI** electrospray ionization

289 **FDR** false discovery rate

290 **FWO** Research Foundation Flanders

291 **GB** gigabytes TODO bits?

292 **HDMS^e** high definition MS^e

293 **HE** high energy

294 **IMS** ion mobility separation

295 **LC** liquid chromatography

296 **LE** low energy

297 **LFQ** label-free quantification

298 **logFC** logarithmic fold change

299 **MRM** multiple reaction monitoring

300 **MS** mass spectrometry

301 **PIM** peptide-ion match

302 **ppm** parts per million

303 **PSM** peptide-spectrum match

304 **ptp** point-to-point

305 **QC** quality control

306 **RAM** random-access memory

307 **RANSAC** random sample consensus

308 **RT** retention time

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

310 **SWIM** single window ion mobility

311 **TOF** time of flight

312 **UGhent** Ghent University

313 **VLAIO** Flanders Innovation & Entrepreneurship

314 **w/w** weight for weight

315 **XIC** extracted ion chromatogram

316 References

- 317 [1] Limonier F, Willems S, Waeterloos G, Sneyers M, Dhaenens M, Deforce D. Estimating the
318 reliability of low-abundant signals and limited replicate measurements through MS2 peak area
319 in SWATH. *PROTEOMICS*;0.ja:1800186.
- 320 [2] Craig R, Beavis RC. TANDEM: matching proteins with tandem mass spectra. *Bioinformatics*
321 2004;20.9:1466–1467. eprint: /oup/backfile/content_public/journal/bioinformatics/
322 20/9/10.1093/bioinformatics/bth092/2/bth092.pdf.
- 323 [3] The M, MacCoss MJ, Noble WS, Käll L. Fast and Accurate Protein False Discovery Rates on
324 Large-Scale Proteomics Data Sets with Percolator 3.0. *Journal of The American Society for*
325 *Mass Spectrometry* 2016;27.11:1719–1727.
- 326 [4] Navarro P, Kuharev J, Gillet LC, Bernhardt OM, MacLean B, Röst HL, Tate SA, Tsou C-C,
327 Reiter L, Distler U, Rosenberger G, Perez-Riverol Y, Nesvizhskii AI, Aebersold R, Tenzer S.
328 A multicenter study benchmarks software tools for label-free proteome quantification. *Nature*
329 *Biotechnology* 2016;34:1130.

1 Material and methods

1.1 Sample preparation

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

1.2 Data acquisition

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

1.3 Peak picking of raw data

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

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

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

1.4 Sample calibration and estimation of ion inter-run differences

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

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

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

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

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

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

1.5 Ion inter-run alignment and noise definition

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

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

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

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

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

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

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

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

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

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

TODO PPM difference calibration between HE-LE

1.7 Aggregate ion network generation

A network was created in which all aggregate ions were vertices. An edge is set between two aggregate ions if and only if they consistently co-elute. Two aggregate ions are defined as *consistently co-eluting* if and only if they co-elute for each overlapping sample. However, as the intra-run differences within each sample are independent, a large sample count can introduce a dimensionality curse, meaning it is unlikely that representative ions co-elute in each sample even if they originate from the same precursor. Therefore the definition of *consistently co-eluting* is weakened to mean that they should have a probability of at least 0.999 to overlap in at least x out of y samples. 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 constraint, two aggregate ions should co-elute in at least two samples to be considered *consistently co-eluting*.

1.8 Fragment-centric annotation of an aggregate ion network

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

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

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

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

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

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

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

505 **1.9 AutoQC, data deposition and source code availability**

506 TODO

507 References

- 508 [1] Limonier F, Willems S, Waeterloos G, Sneyers M, Dhaenens M, Deforce D. Estimating the
509 reliability of low-abundant signals and limited replicate measurements through MS2 peak area
510 in SWATH. *PROTEOMICS*;0.ja:1800186.
- 511 [2] Craig R, Beavis RC. TANDEM: matching proteins with tandem mass spectra. *Bioinformatics*
512 2004;20.9:1466–1467. eprint: /oup/backfile/content_public/journal/bioinformatics/
513 20/9/10.1093/bioinformatics/bth092/2/bth092.pdf.
- 514 [3] The M, MacCoss MJ, Noble WS, Käll L. Fast and Accurate Protein False Discovery Rates on
515 Large-Scale Proteomics Data Sets with Percolator 3.0. *Journal of The American Society for*
516 *Mass Spectrometry* 2016;27.11:1719–1727.
- 517 [4] Navarro P, Kuharev J, Gillet LC, Bernhardt OM, MacLean B, Röst HL, Tate SA, Tsou C-C,
518 Reiter L, Distler U, Rosenberger G, Perez-Riverol Y, Nesvizhskii AI, Aebersold R, Tenzer S.
519 A multicenter study benchmarks software tools for label-free proteome quantification. *Nature*
520 *Biotechnology* 2016;34:1130.

521 **Supplementary**

522 **Peak picking**

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

Pseudo aggregate formation

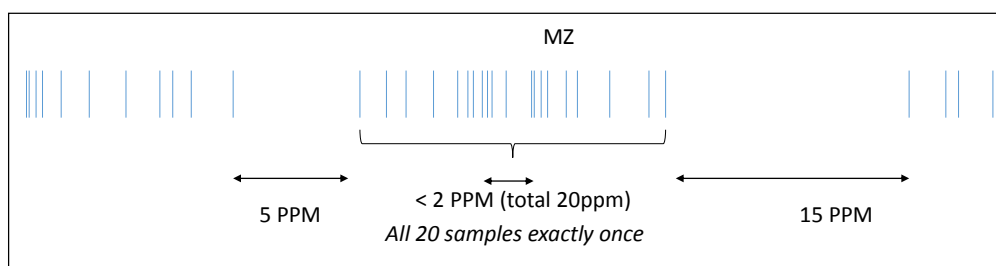


Figure 4: Schematic overview of quick alignment to detect pseudo aggregate ions.

524 **Schematic overview of pseudo aggregate ion group piece-wise linear transformations**
525 **for retention times**

RT Calibration

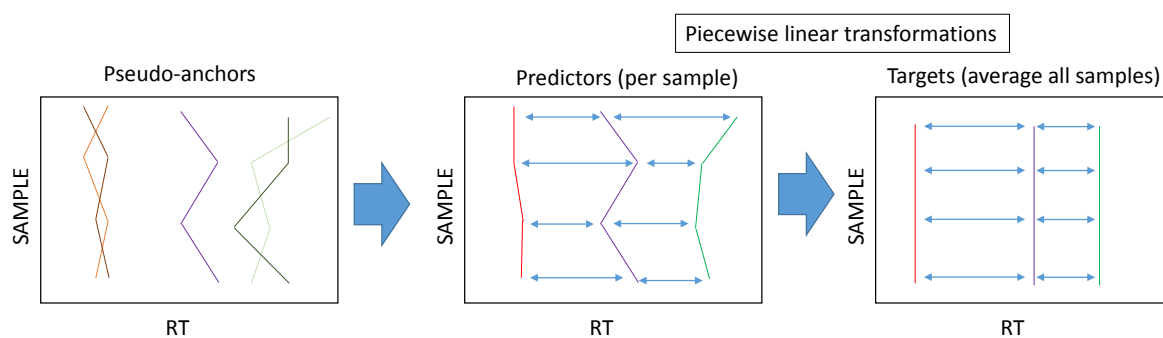


Figure 5: Schematic overview of pseudo aggregate ion group piece-wise linear transformations

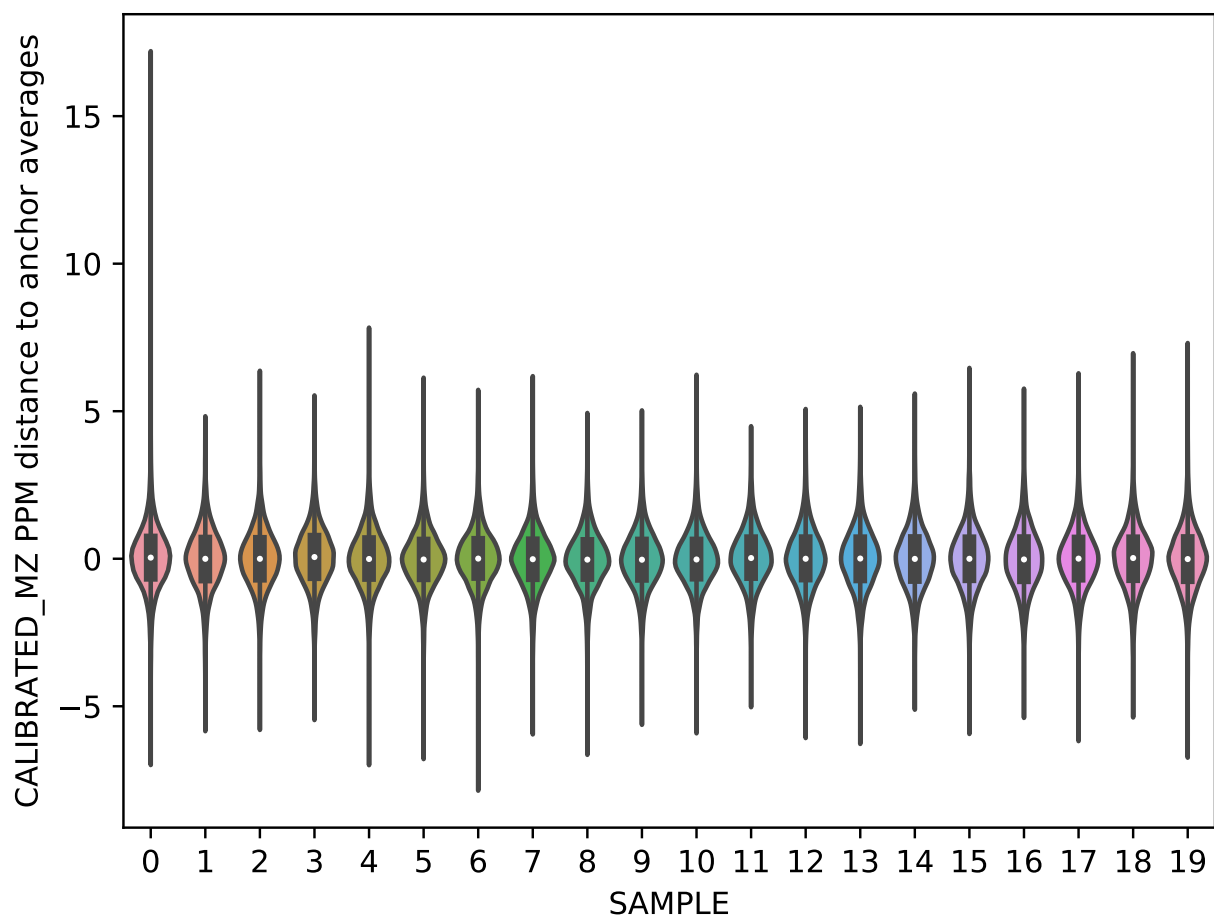


Figure 6: Estimation of m/z error per sample toward the pseudo aggregate ions.

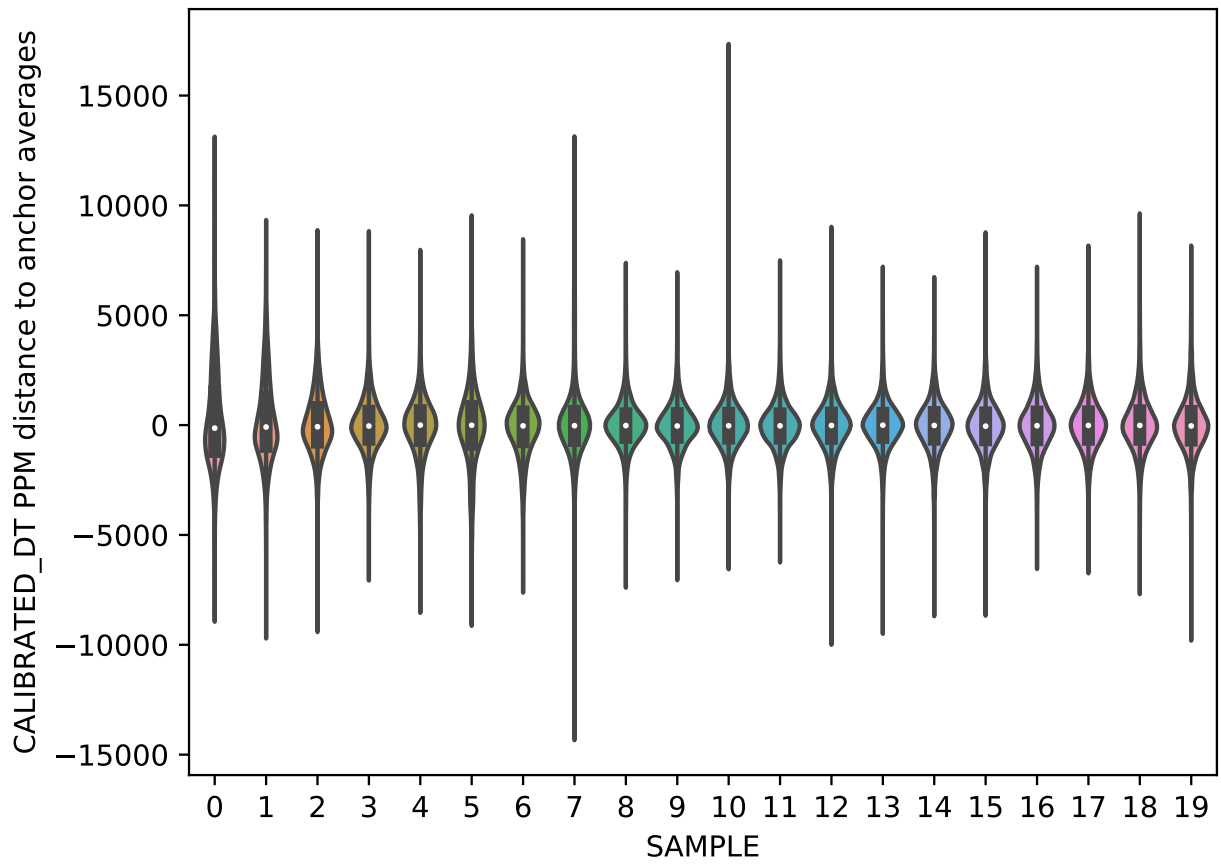


Figure 7: Estimation of DT error per sample toward the pseudo aggregate ions.

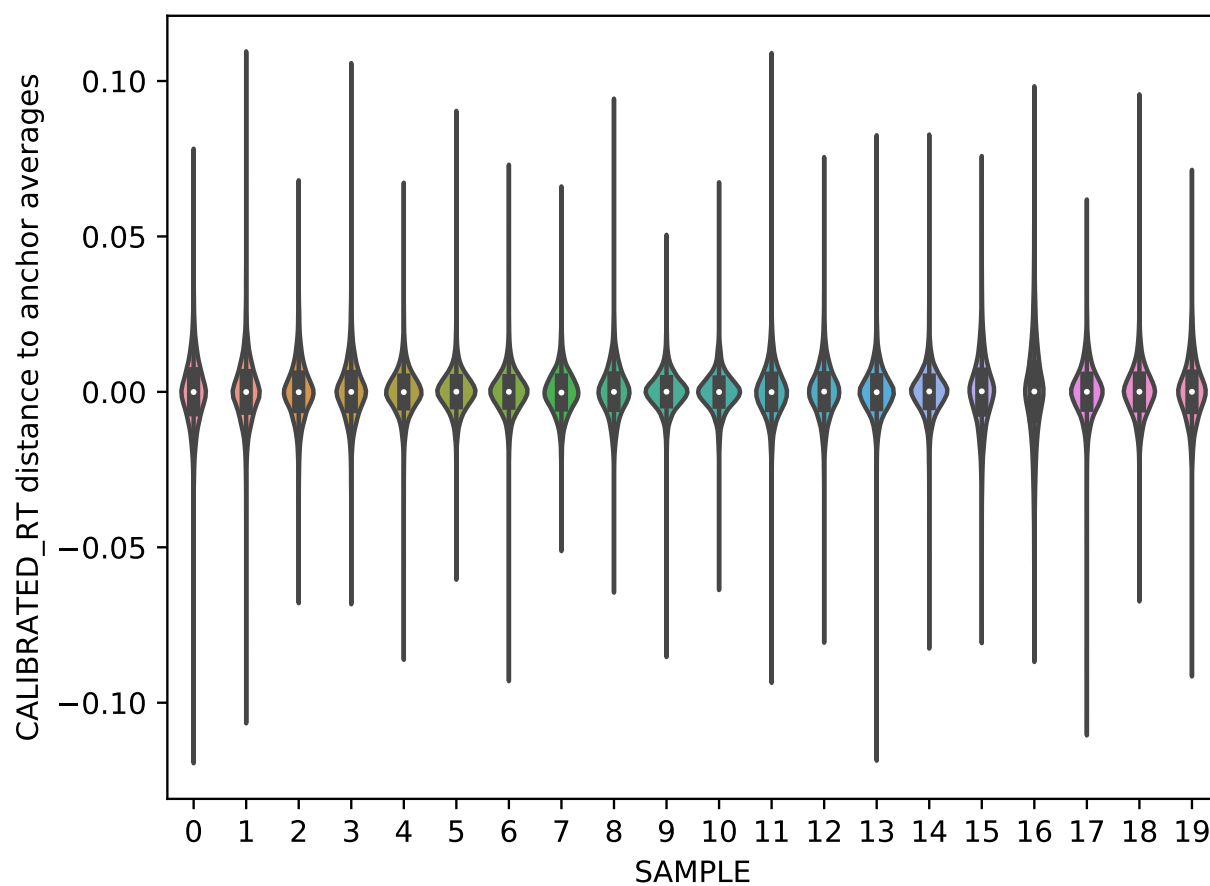


Figure 8: Estimation of RT error per sample toward the pseudo aggregate ions.

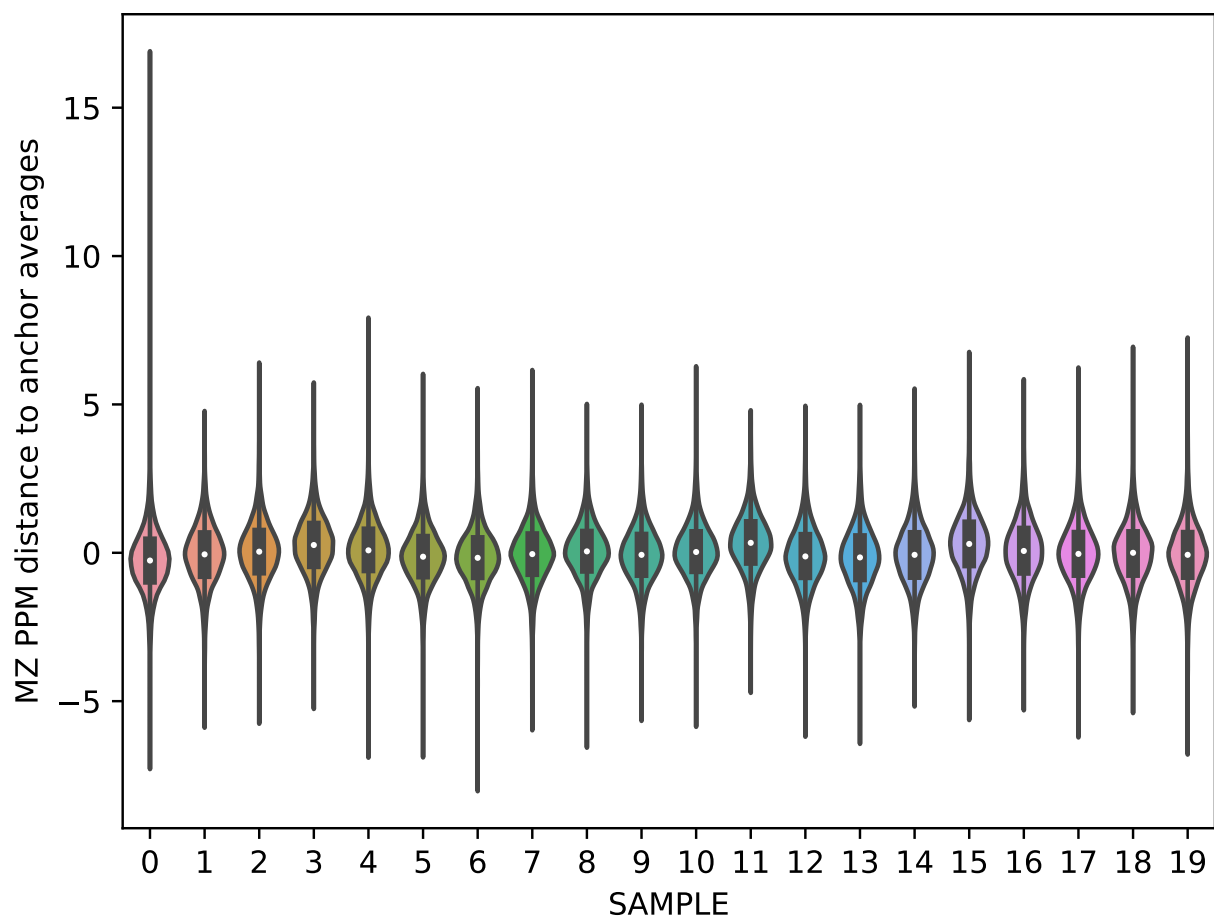


Figure 9: Uncalibrated m/z error per sample toward the pseudo aggregate ions.

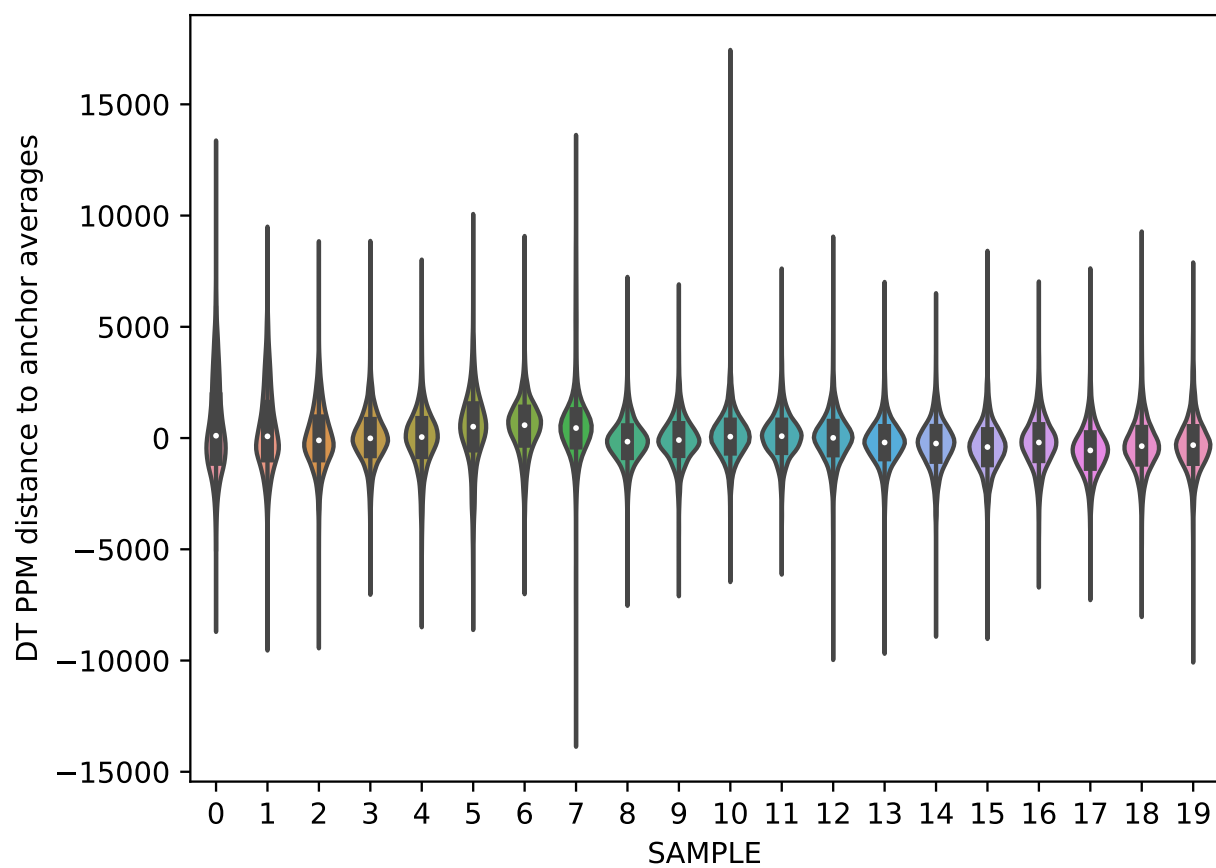


Figure 10: Uncalibrated DT error per sample toward the pseudo aggregate ions.

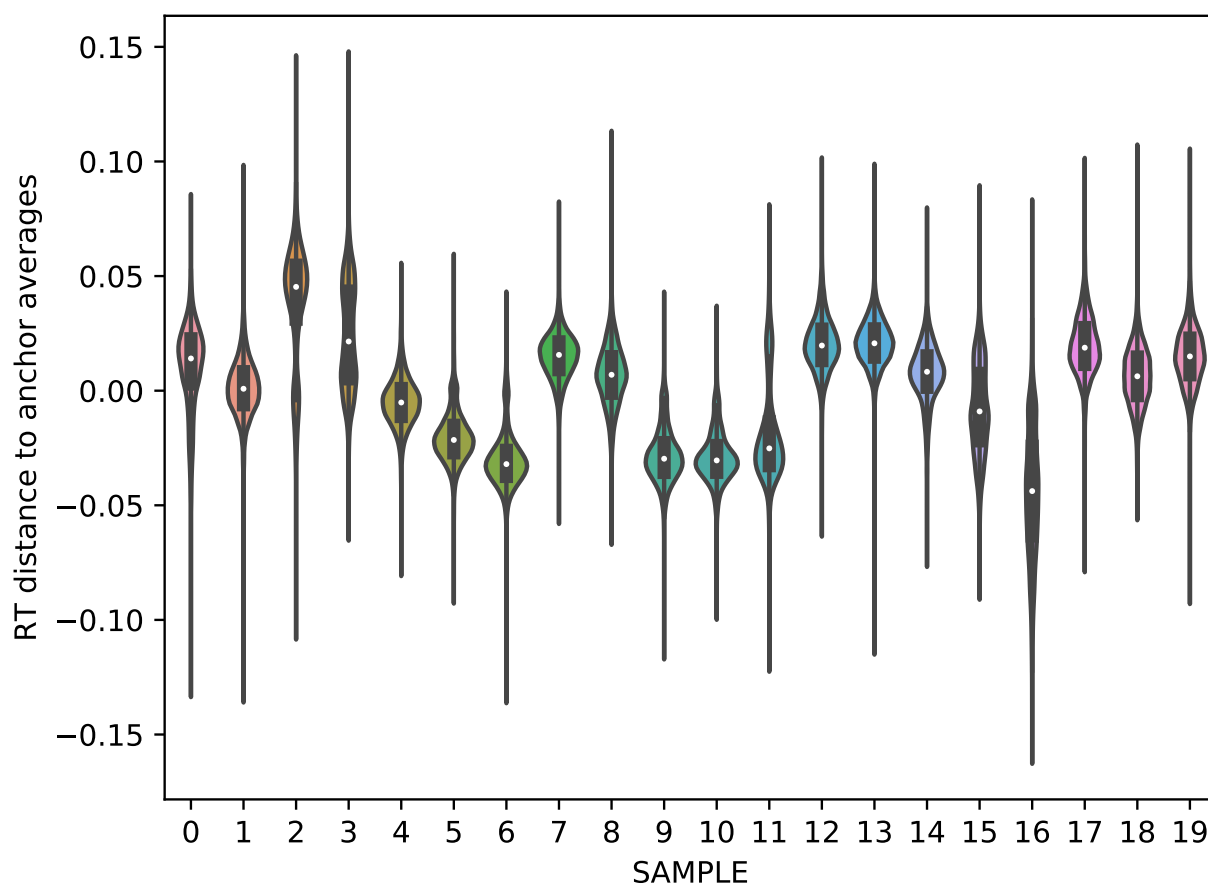


Figure 11: Uncalibrated RT error per sample toward the pseudo aggregate ions.

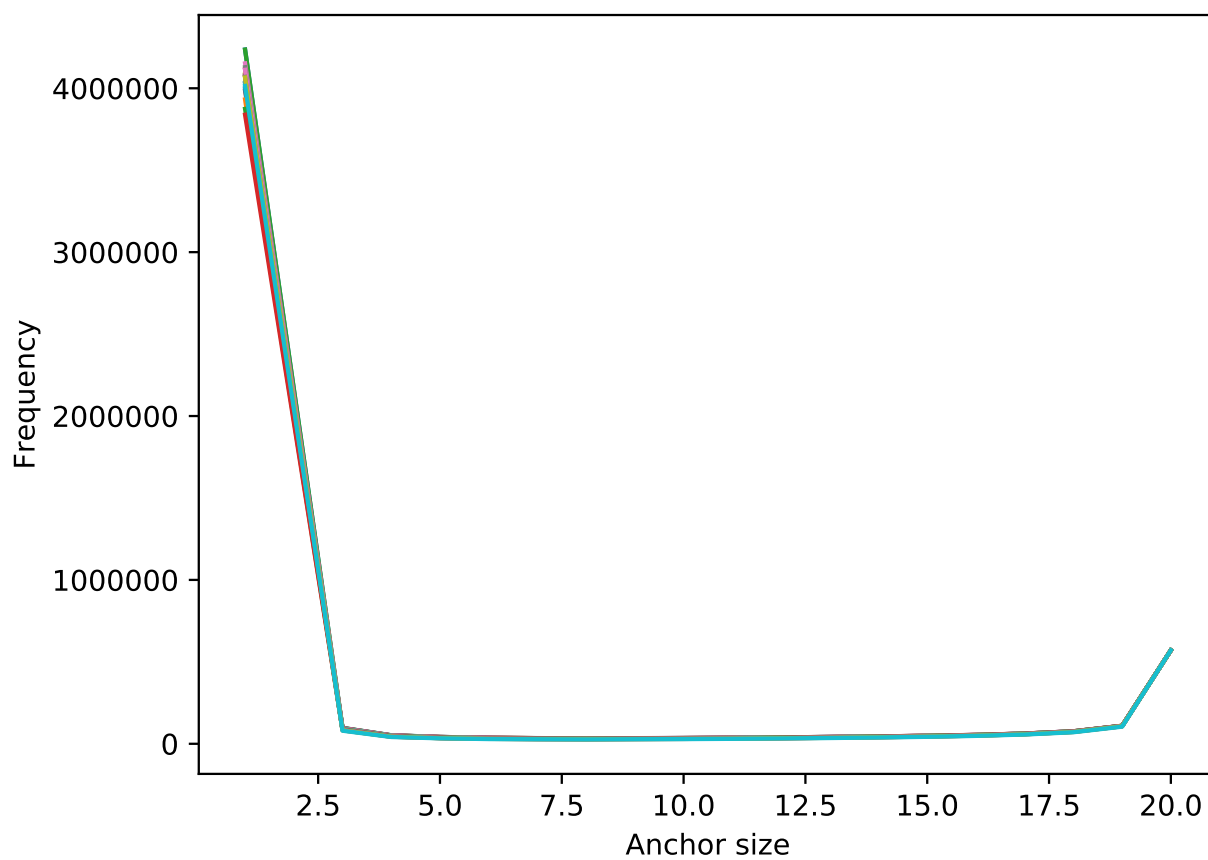


Figure 12: Low energy aggregate ion sizes.

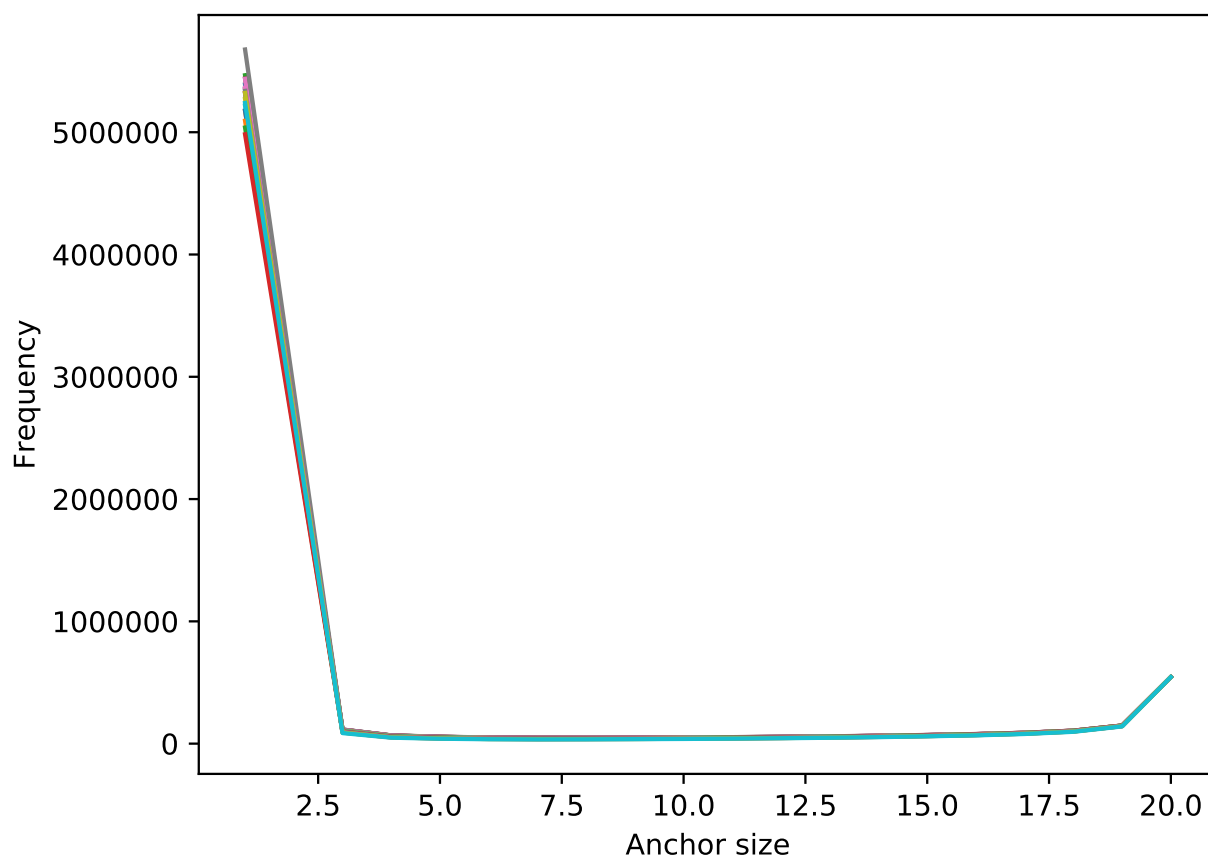


Figure 13: High energy aggregate ion sizes.

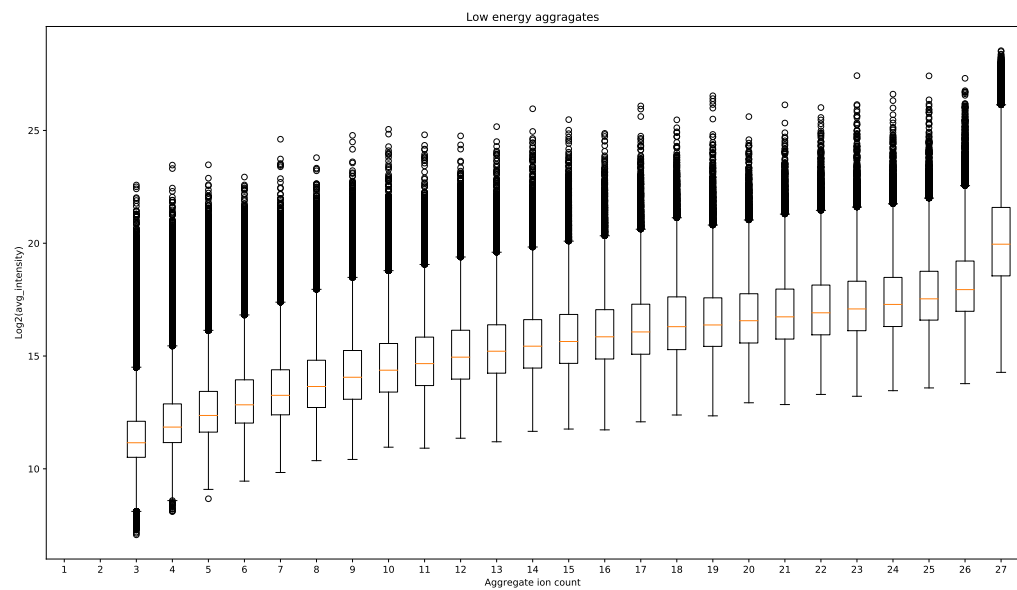


Figure 14: Low energy aggregate ion intensities.

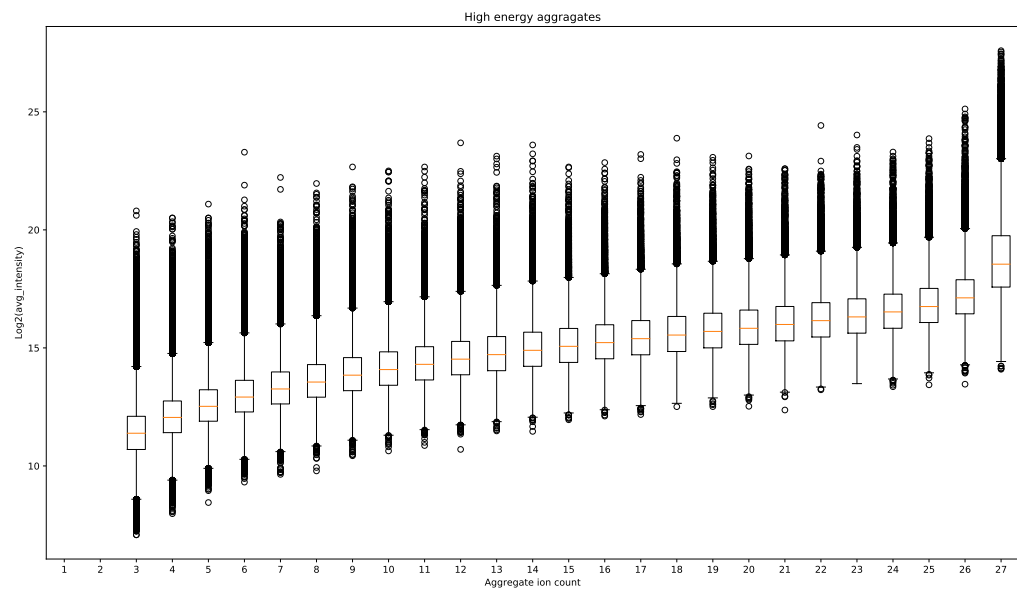


Figure 15: High energy aggregate ion intensities.

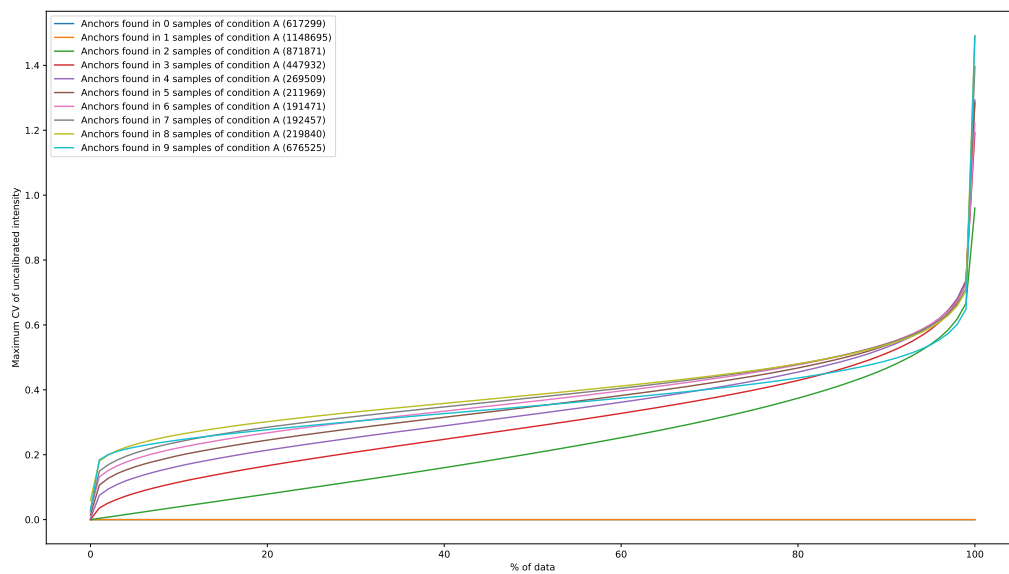


Figure 16: CV of uncalibrated intensities in condition A.

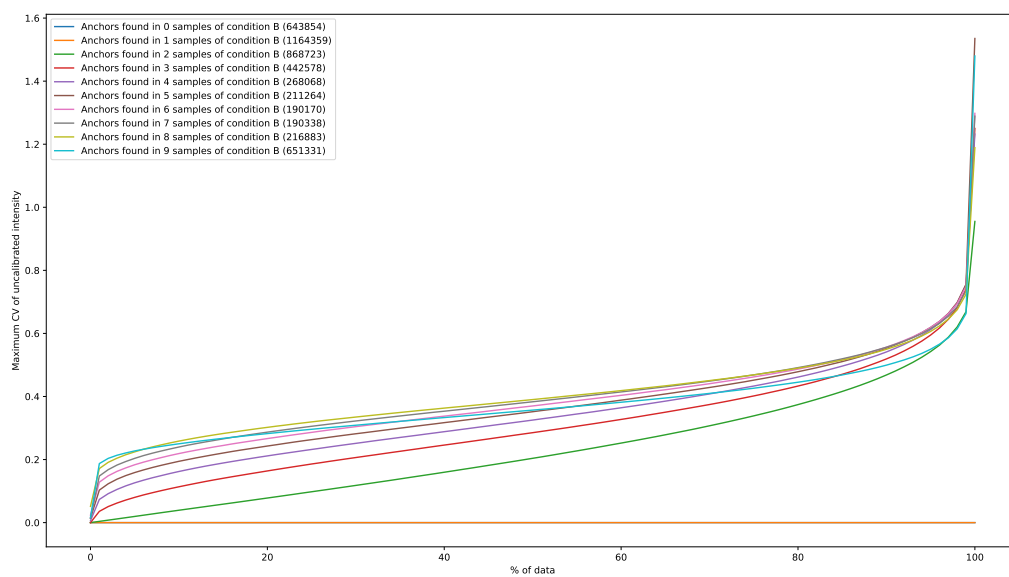


Figure 17: CV of uncalibrated intensities in condition B.

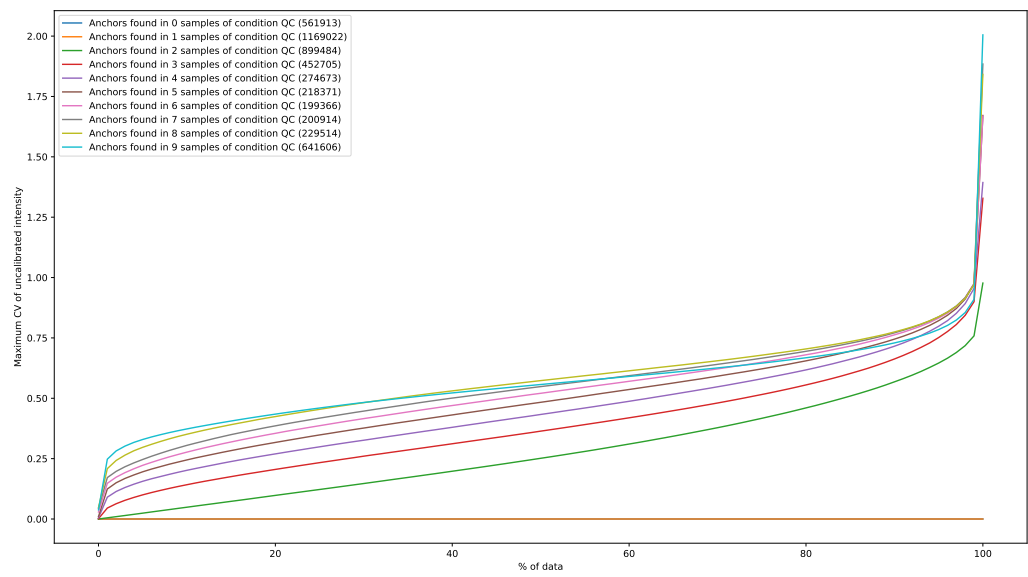


Figure 18: CV of uncalibrated intensities in condition QC.

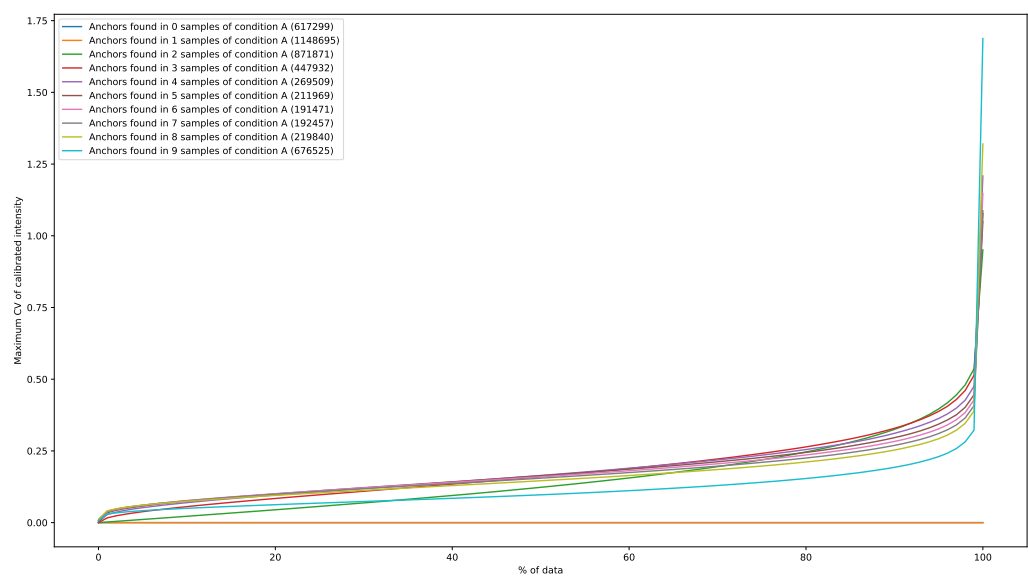


Figure 19: CV of calibrated intensities in condition A.

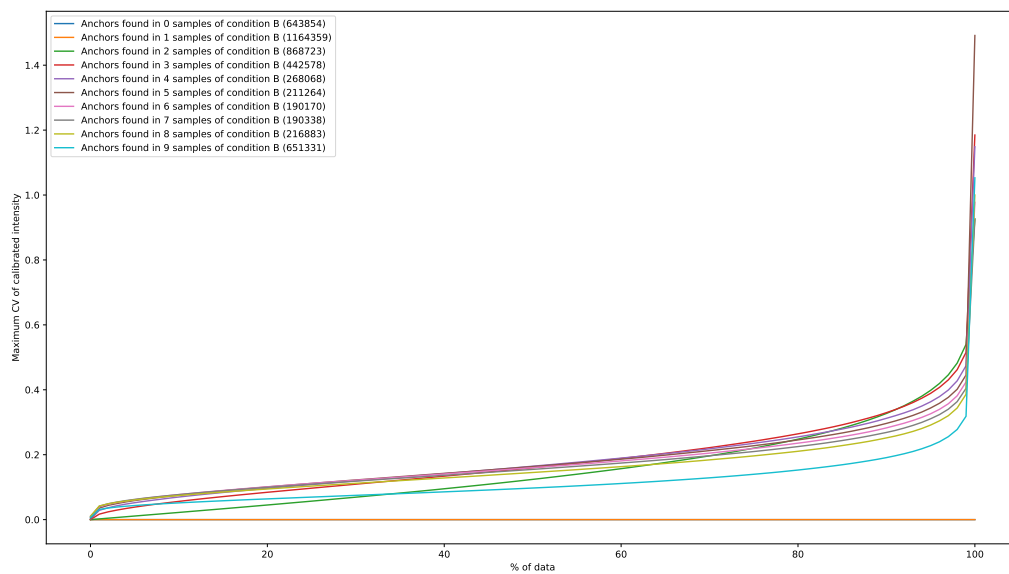


Figure 20: CV of calibrated intensities in condition B.

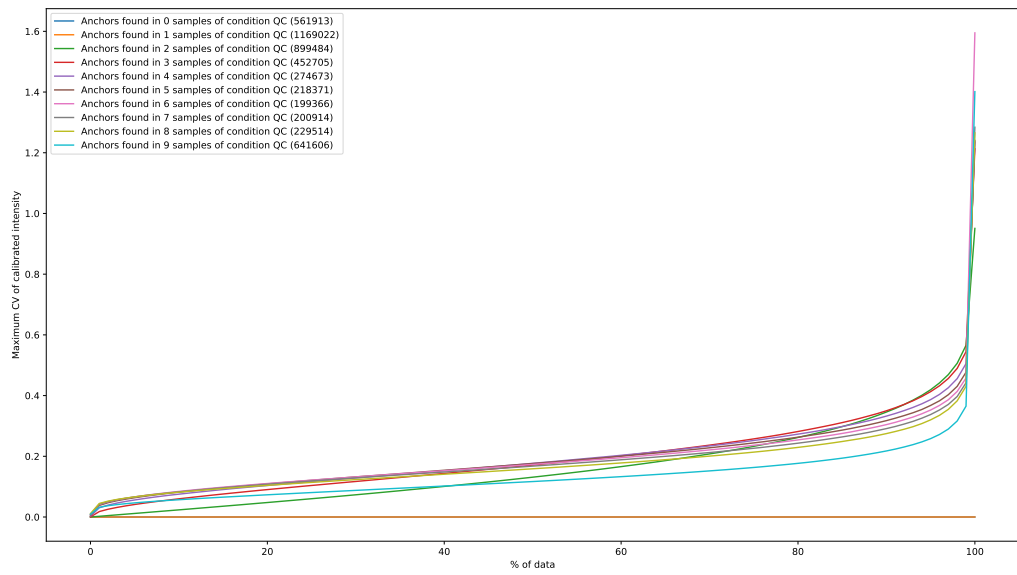


Figure 21: CV of calibrated intensities in condition QC.

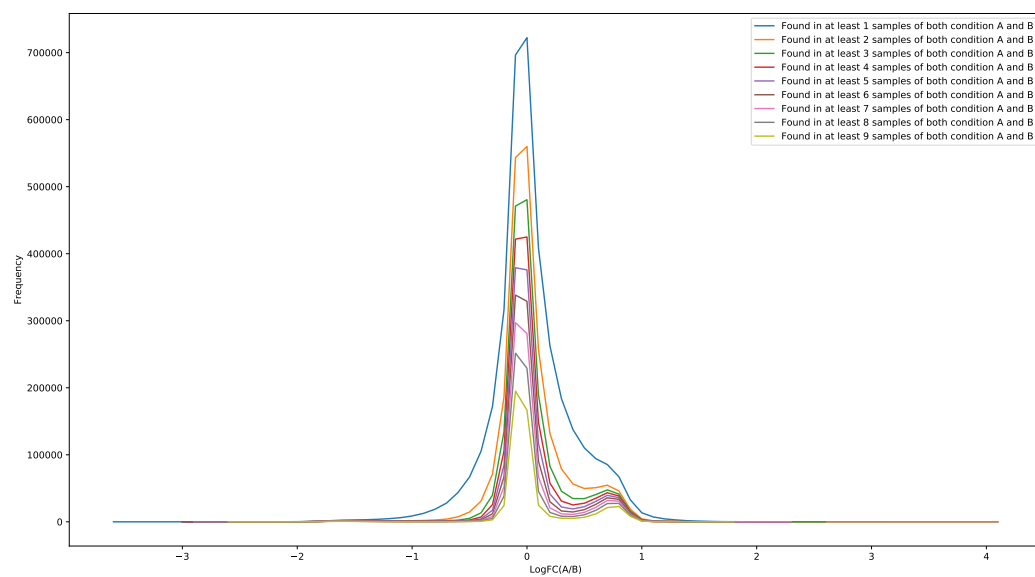
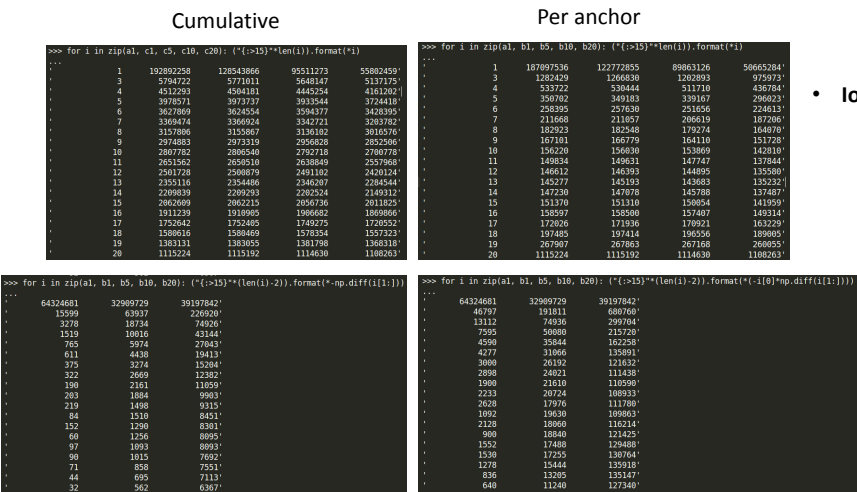


Figure 22: Distribution of logarithmic fold changes.

Signal detection



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

Figure 23: HistoPyA’s robustness to peak picking noise.

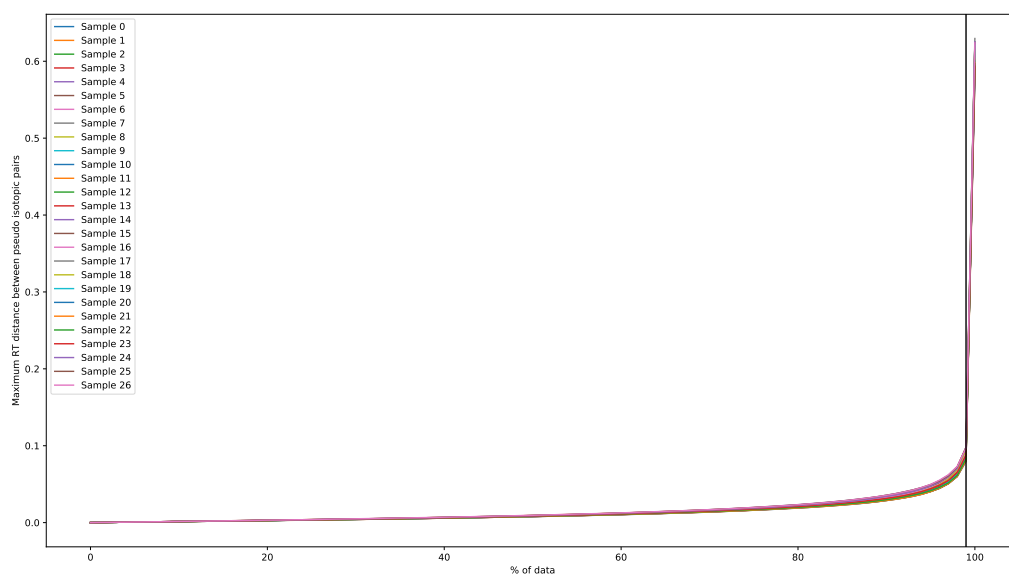


Figure 24: Distribution of intra-run retention time distances.

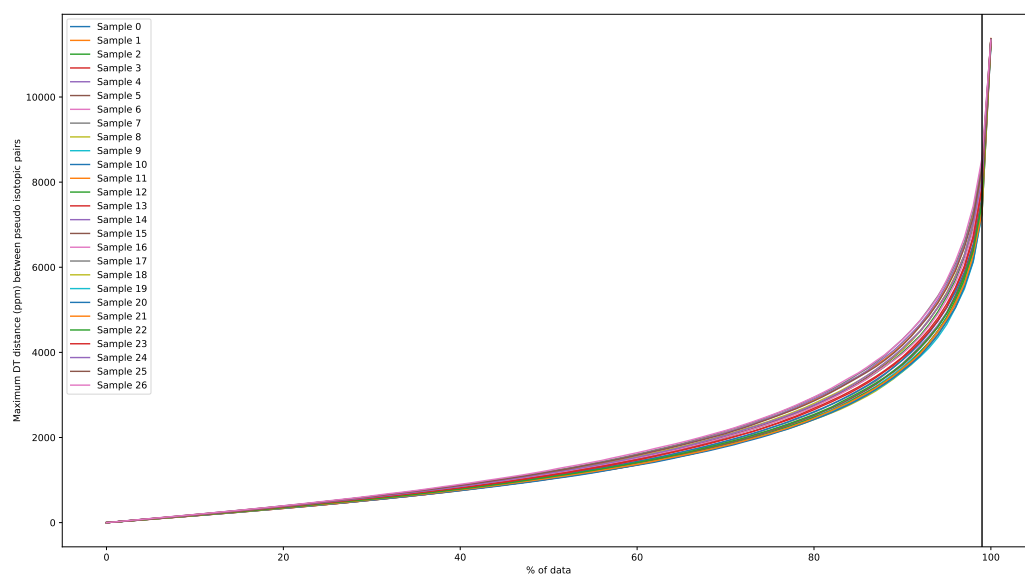


Figure 25: Distribution of intra-run drift time distances.

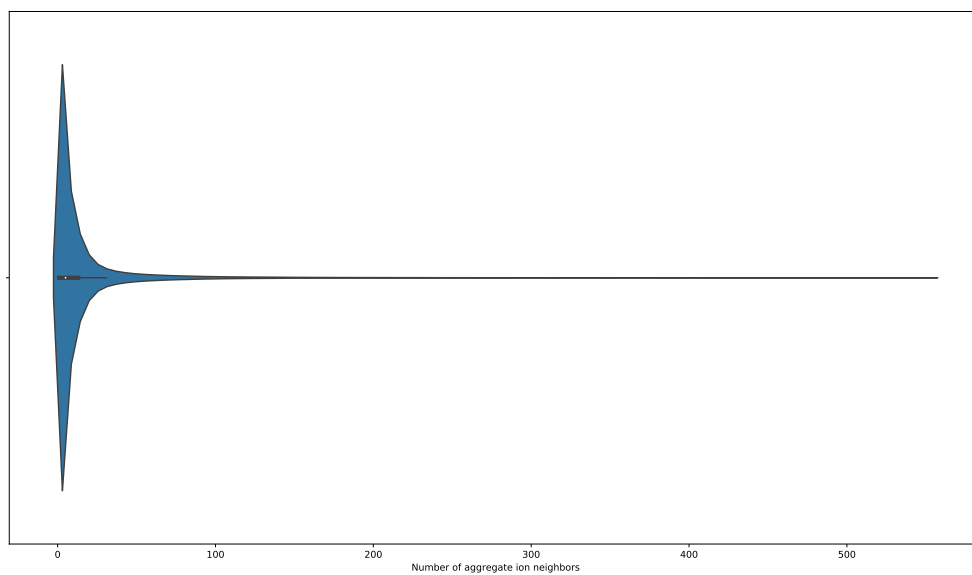


Figure 26: Distribution of aggregate ion neighbor counts.



Figure 27: Sample scheme.



Figure 28: Drift time shift between low energy and high energy channels.