1 **Title of Paper**

2 **DIAlignR provides precise retention time alignment across distant runs in DIA and**

3 **targeted proteomics**

4

5 **Authors**

6 Shubham Gupta[1,2], Sara Ahadi[3], Wenyu Zhou[3], Hannes Röst[1,2*]

7

8 **Affiliations:**

9 [1]Department of Molecular Genetics, University of Toronto, Toronto, ON M5G 1A8, Canada

10 [2]Donnelly Centre for Cellular and Biomolecular Research, University of Toronto, Toronto, ON

11 M5S 3E1, Canada

12 [3]Department of Genetics, Stanford University School of Medicine, Stanford, CA 94305, USA.

13 *Corresponding author

14 Hannes Röst, 160 College Street, Toronto, ON, M5S 3E1, Canada

15 Email: hannes.rost@utoronto.ca

16 **Abbreviations:**

17 AUC       Area Under the Curve

18 DIA        Data-independent acquisition

19 LC         Liquid chromatography

20 LOESS   Local weighted regression

21 RSE       Residual Standard Error

22 RT         Retention time

23 XIC        Extracted ion chromatograms

24

25 **Running Title:**

26 DIAlignR: Mapping retention time using MS2 chromatograms

1

1

## 2    **Data Availability:**

3    Raw chromatograms and features extracted by OpenSWATH are available on PeptideAtlas.

4    Servername: ftp.peptideatlas.org

5    Username: PASS01280

6    Password: KQ2592b

**1    Abstract**

2    SWATH-MS has been widely used for proteomics analysis given its high-throughput and

3    reproducibility but ensuring consistent quantification of analytes across large-scale studies of

4    heterogeneous samples such as human-plasma remains challenging. Heterogeneity in large-

5    scale studies can be caused by large time intervals between data-acquisition, acquisition by

6    different operators or instruments, intermittent repair or replacement of parts, such as the liquid

7    chromatography column, all of which affect retention time (RT) reproducibility and successively

8    performance of SWATH-MS data analysis. Here, we present a novel algorithm for retention time

9    alignment of SWATH-MS data based on direct alignment of raw MS2 chromatograms using a

10    hybrid dynamic programming approach. The algorithm does not impose a chronological order of

11    elution and allows for alignment of elution-order swapped peaks. Furthermore, allowing RT-

12    mapping in a certain window around coarse global fit makes it robust against noise. On a

13    manually validated dataset, this strategy outperforms the current state-of-the-art approaches. In

14    addition, on a real-world clinical data, our approach outperforms global alignment methods by

15    mapping 98% of peaks compared to 67% cumulatively and DIAlignR can reduce alignment error

16    up to 30-fold for extremely distant runs. The robustness of technical parameters used in this

17    pairwise alignment strategy has also been demonstrated. The source code is released under

18    the BSD license at https://github.com/Roestlab/DIAlignR.

19

**20    Introduction:**

21         In translational research, protein biomarkers and therapeutic targets are usually

22    discovered by data-driven methods such as by linking protein abundance patterns with disease

23    conditions. A large sample cohort is essential in these studies as substantial biological variability

24    exists in the population and enough statistical power is required to identify disease specific

25    events[1,2]. Blood plasma is a good source of clinical information of a patient as it can be obtained

26    noninvasively and proteins from affected tissue can potentially leak into the blood. Plasma

1 samples, unfortunately, are highly challenging for proteomic analysis due to the diversity of

2 peptides within the samples and high dynamic range of plasma proteins[3]. Therefore,

3 quantification of plasma proteins requires a highly reproducible reduction of complexity and

4 measurement within a wide dynamic range. The situation is exacerbated across large-scale

5 studies which makes development of plasma biomarker challenging[2,3].

6       In the past two decades, mass spectrometry (MS) based proteomics has made rapid

7 advances and high degree of innovation in obtaining identification and quantification of proteins

8 in various biological samples[2,4]. Targeted proteomics methods, specifically selected reaction

9 monitoring (SRM), can provide high reproducibility across multiple runs. However, it is limited by

10 low throughput and can measure abundance of only a few tens to low hundreds of proteins per

11 study[1,5].

12       Recently, we developed SWATH-MS, an approach for targeted analysis of data-

13 independent acquisition (DIA) data, which can reproducibly quantify larger sets of peptides in

14 large-scale clinical studies[5,6]. Implementing this method in the clinical field could provide

15 comprehensive characterization of sample across various conditions. It has allowed to

16 reproducibly quantify about 2000 proteins in a biomarker study on tumorous kidney and healthy

17 tissues[1,7] and has the potential to record a molecular inventory of samples comprising a large

18 number of proteotypes, thus making longitudinal monitoring of a patient possible[1].

19       In DIA mode, precursors in MS1 are selected for a predetermined m/z range and

20 fragmented non-specifically. This produces multiplexed MS2 spectra of fragment-ions of all

21 selected precursors. The DIA data can be analyzed by using either a library-based approach[5,8]

22 or a library-free approach[9]. Library-based approaches have shown to be capable of accurate

23 peptide and protein quantification in complex samples[5,10,11]. Nonetheless, obtaining reproducible

24 and robust analysis of clinical plasma samples has been challenging even with SWATH-MS, as

25 large variations in number of proteins in individual runs were observed[5,10,11]. One of the major

26 factors driving variability is the retention time deviation between assay library and plasma

4

1    peptides' elution profiles. In experiments carried out by Nigjeh and coworkers, most of the

2    peptides had RT variation of about 10 minutes between technical replicates, affecting the

3    robustness of peptide quantification[3]. This variation, if left uncorrected, may also result into

4    incorrect and inconsistent identification of the peptides[3].

5         Current DIA data analysis software use iRT peptides to calculate a monotonic retention

6    time function (linear regression[12] or segmented regression[13]) with respect to a library. Using this

7    mapping, extracted ion chromatograms (XICs) from MS2 spectra are obtained for peak-picking.

8    Software usually finds multiple potential peak-groups in XICs, which makes downstream

9    analysis challenging. By establishing peak correspondence among runs, correct peptide elution

10   time could be determined for each MS run[5,14].  A shift in retention time (RT) is often considered

11   as a system-level variation which is modelled using monotonic functions between two runs[14].

12   However, this assumption may not always be accurate, and specifically among distant runs,

13   singularities specific to a single peptide are common that produces relative peak switching

14   where the elution order of two peptides is swapped across two runs[14–16]. This phenomenon is

15   increasingly likely in larger studies and very probable in large-scale clinical studies in which

16   data-acquisition happens over a span of years.

17        There are many methods in the literature for establishing correspondence in retention

18   times. Current RT alignment algorithms in metabolomics and proteomics were mostly developed

19   before the development of SWATH-MS[14] and, therefore, rely on either MS1 chromatograms[17–23],

20   picked features[24–27], or a combination of both [28,29]. These algorithms usually find a global

21   pairwise alignment function using dynamic programming on raw MS1 chromatograms[20–23,28] or

22   on feature-lists[24,25] which include methods using band constraints. For complex samples, so-

23   called "landmark peaks"[28,29] have been used to improve RT alignment accuracy. However, most

24   of these approaches rely on MS1 data and the resulting alignment functions are influenced by

25   all constituent peptides. In SWATH-MS runs, MS2 data has high signal-to-noise ratio and is

26   reproducible across multiple runs. Previous research on RT alignment of SWATH runs has

1    relied on MS2 feature finding software. These either align MS2 features using bipartite

2    matching[16] or use the features to calculate a global function by local weighted regression

3    (LOESS)[30] or by a kernel density approach[5,31] (see Supplemental Section S1). These

4    approaches, however, provide suboptimal results in case of high noise, missing features or

5    when feature detection algorithms malfunction. Furthermore, the global monotone functions do

6    not account for peptide switching as a monotone function disallows retention time reversal

7    between any two peptides[16].

8         Here we present DIAlignR, a retention time alignment algorithm that addresses these

9    shortcomings of previous methods. Our algorithm does not require features and is capable of

10   directly aligning the raw multiplexed MS2 chromatographic traces from targeted proteomics

11   data. Our approach uses dynamic programming to obtain an optimal mapping between

12   chromatograms which contain local information as multiple, close-by peaks around the eluted

13   peak-group. Independent RT alignment of each precursor facilitates the alignment of elution-

14   order swapped peaks. Our method is also capable of using a global whole-run alignment for

15   guidance, making it robust against noise. Thus, DIAlignR can flexibly handles user preference of

16   selecting between extremes of global and local alignment.

17        We provide free-access to our source-code and our R-package at

18   https://github.com/Roestlab/DIAlignR. We have tested our algorithm on a manually validated

19   dataset of over 7000 chromatograms and demonstrate improved performance over existing

20   methods. We have also tested our algorithm on 24 randomly selected blood plasma runs,

21   selected from a heterogeneous cohort measured across many months. For both datasets, our

22   algorithm outperforms global alignment methods and is capable of correcting mis-annotations

23   introduced by feature detection algorithms. For very distant runs, it could also precisely align

24   switched peaks which is not possible using global alignment methods[8,12].

25

26   **Material and Methods**:

1   Alignment algorithms are useful for mapping signals, which are close to noise, over

2   several runs. Few datasets exist that can be used for benchmarking as often the ground

3   truth is unknown.

4   <u>Validation dataset</u>

5   For benchmarking, we have used a previously published dataset[8] of 16 SWATH runs

6   from *Streptococcus pyogenes* bacterial strains. In these runs, 452 randomly selected precursors

7   were manually annotated[5]. Out of these, eight precursors have annotation for less than two

8   runs, and were thus removed. Seven other precursors from the remaining set have annotated

9   peaks outside of the XICs, making them inapplicable and, hence, were removed from

10  benchmarking (see the supplemental Section S2). Therefore, 437 annotated precursors are

11  considered for performance testing of the developed DIAlignR tool against global alignment

12  approaches. Annotated retention times of the precursors are available in the Supplemental

13  Table 1. Since, annotated peaks were selected randomly from the validation dataset, this

14  dataset has 4.9% peaks with signal-to-noise ratio less than 1 (Supplemental Figure 1*a*).

15

16  <u>Large-scale human plasma dataset</u>

17  We have performed SWATH-MS on 975 human plasma samples in 12 batches

18  from 17 February 2017 to 20 July 2017 (IRB 23602). Tryptic peptides of plasma samples were

19  separated on a NanoLC 425 System (SCIEX). 5ul/min flow was used with trap-elute setting

20  using a 0.5 x 10 mm ChromXP (SCIEX). LC gradient was set to a 43-minute gradient from 4-

21  32% B with 1-hour total run.  Mobile phase A was 100% water with 0.1% formic acid. Mobile

22  phase B was 100% acetonitrile with 0.1% formic acid. 8ug load of undepleted plasma on 15cm

23  ChromXP column. MS analysis was performed using SWATH Acquisition on a TripleTOF 6600

24  System equipped with a DuoSpray Source and 25μm I.D. electrode (SCIEX). Variable Q1

25  window SWATH Acquisition methods (100 windows) were built in high sensitivity MS/MS mode

26  with Analyst TF Software 1.7.

1     To reduce the number of pairwise alignment, randomly two runs from each batch were

2     selected; their metadata and OpenSWATH output files are described in the supplemental Table

3     2. Since peaks were not visually validated, only peaks with low FDR score were considered for

4     performance evaluation. Therefore, peaks of target precursors with a q-value less than $10^{-3}$ (m-

5     score < 1e-03, peak-group rank = 1) were selected and precursors were required to be present

6     in all 24 runs. Successively, fragment-ion chromatograms of selected 406 precursors were

7     extracted and parsed using OpenSWATH[8,32] and "mzR" package[33] with default parameters. The

8     retention time of the peptides in all 24 runs is provided in the Supplemental Table 3. The

9     chromatograms are available on PeptideAtlas (ftp.peptideatlas.org PASS01280:KQ2592b).

10     A tabular description of both datasets is provided as Table 1.

11

12     **<u>Chromatogram Alignment Algorithm</u>**

13     In targeted proteomics or SWATH-MS experiments, each precursor is measured using

14     one or more fragment ions (transitions). In general, we recommend for DIA / SWATH-MS

15     analysis to use at least six fragment ions[4]. For each fragment-ion an extracted-ion

16     chromatogram (XIC or chromatogram) is obtained. A collection of one or more chromatograms

17     is called a "chromatogram group" which maps to the given precursor. If a precursor is measured

18     using $n$ transitions, then for each run the respective "chromatogram group" consists of $n$ XICs,

19     which is the raw data for our alignment procedure.

20     A chromatogram group can be considered a collection of time-series signals. The

21     similarity of the time-series signals between chromatogram groups from runA (ChromA) and

22     runB (ChromB) can be calculated. If a precursor has $n$ fragment-ions, and each XIC has $I$ and $J$

23     time-points in *ChromA* and *ChromB*, respectively as shown in Fig. 1*a*, the similarity between all

24     time-points can be represented as a similarity matrix $s$ (Fig. 1*b* and *c*). Thus,

25     $$s = f(ChromA, ChromB).$$

26     The function $f$ is termed as a similarity measure and can be selected by the user (see below).

1    **a.  Similarity measure:**

2        In our R package DIAlignR, we have implemented several similarity measures which

3    have been suggested in previous literature for chromatograms such as covariance, dot-product,

4    Pearson's correlation, spectral angle and euclidean distance[8,28]. We have observed that the dot-

5    product between all *I* and *J* data-points provide information about both magnitude and angle

6    between two time-points, hence segregating elution signal from the background. If each data

7    point of chromatogram is represented by a vector in $n$ dimensional space ($n = 3$ in Fig. 1*a*), the

8    resulting dot-product of the two vectors will be as shown in Fig. 1*b*. Thus, with dot-product

9    similarity, the matrix *s* from all vectors of both chromatogram-groups is defined as,

10
$$s_{ij} = \sum_{k=1}^{n} a_{ik} b_{jk}$$

11    Where $i \in \{1, \ldots, I\}$ and $j \in \{1, \ldots, J\}$ represents index of vectors in *ChromA* and *ChromB*,

12    respectively. A color-coded similarity matrix of size *I x J* is shown in Fig. 1*c*. However, to reduce

13    the impact of noise peaks, a modified dot-product, termed as "Masked dot-product", is used

14    where higher similarity scores are checked again for spectral angle similarity (see the

15    Supplemental Section S5). A path in the resulting similarity matrix is calculated using dynamic

16    programming which directly translates to a retention time alignment that maps indices/time from

17    *ChromA* to *ChromB* and vice-versa.

18

19    **b.  Penalizing similarity matrix with global alignment:**

20        While dynamic programming will find a path of the highest cumulative score, in some

21    instances the score is driven by alignment to noise and can lead to a solution where the

22    alignment is highly divergent from a global linear or non-linear alignment. To make the

23    alignment robust against noise and in order to incorporate information from the global context,

24    we have added an option in our algorithm to modify the similarity matrix *s* (Fig. 1*d*) using

25    feature-based global alignment such as LOESS. Residual standard error of the fit is utilized to

1  define a region of non-interference in the similarity matrix and values outside of it punished with

2  negative score (see the Supplemental Section S5). This allows us to find an alignment path

3  within a reasonable time window relative to global prediction and avoid large deviations.

4

5  **c. Overlap Alignment with affine gap penalty:**

6  The optimal alignment path is found by recursively calculating all possible optimal paths

7  from the start of the similarity matrix (1,1) to the end of it (I, J) using dynamic programming[34].

8  Chromatogram-groups *ChromA* and *ChromB* may not have end-to-end mapping as these could

9  be partial chromatograms which were extracted around the expected peptide elution (as

10  determined by iRT peptides for example). Therefore, overlap alignment instead of a global

11  alignment of MS2 chromatogram groups is employed. This approach allows free end-gaps and

12  thus, allows to slide chromatograms freely without incurring any gap-penalty for it.

13  To widen or shrink chromatogram peaks, a gap of unit length is a reasonable choice as it

14  will distribute gaps along the complete peak. Therefore, an affine gap penalty scheme is utilized

15  with higher gap penalty for gap length of more than one. In this approach, three matrices (Matrix

16  M, A, B) are defined which recursively calculates score for gaps of more than unit length[34]. The

17  overlap alignment path using affine gap-penalty is presented in Fig. 1*e*. The running time of

18  such alignment is O($max(I, J)^3$). A heuristic data-driven approach is employed to obtain suitable

19  affine gap penalties from the similarity matrix (see the Supplemental Section S5). Mapping the

20  alignment path to the initial time values provides aligned chromatograms as depicted in Fig. 1*f*.

21

22  **d. Running time for alignment:**

23  Alignment of MS2 chromatograms of each peptide/precursor has running time of order

24  O(max(I, J)^3); however, chromatograms of different precursors can be aligned independently.

25  Therefore, we employ parallelization for different peptides to obtain much faster speed for

26  complete run time-mapping.

1

2      **e.  Optimization of algorithm parameters:**

3      There are various parameters used in DIAlignR. A description of these parameters is

4  available in the Supplemental Section S5. We have employed validation dataset for parameter

5  optimization and have used the number of peaks aligned within half chromatographic peak-

6  width and cumulative RT alignment error as optimization target.

7

8  **Performance metrics for comparison with current algorithms:**

9      We have used the manually validated dataset[5] to compare DIAlignR to the current state-

10  of-the-art method (e.g. TRIC[5]) which utilizes a set of high confidence peaks ("anchor peptides")

11  to compute a linear or non-linear alignment function that transforms RT values from run1 to

12  run2. We have chosen LOESS (local regression) as well as linear regression for evaluation. For

13  LOESS, both optimized spanvalue from cross-validation (as used in TRIC) and default

14  spanvalue (= 0.75) of the R software environment are tested[30]. For LOESS fit, ⅓ cross-

15  validation is performed to obtain the optimum span value between two runs[5,30]. Steps to obtain a

16  global fit (monotone mapping function) are detailed in the Supplemental Section S3.

17      Retention time error is calculated by comparing against the manual annotation of the *S.*

18  *pyogenes* dataset[5] and the resulting distribution of the number of peptides aligned within a

19  certain RT tolerance is used as a measure of overall accuracy of the alignment algorithm.

20  Manual annotations are not available for the human plasma dataset, therefore, the high-quality

21  results (peaks with low FDR cutoff) of OpenSWATH are used for benchmarking.

22

23  **Results**:

24  *Parameter optimization.*

25      Here, we present an algorithm for multi-trace chromatographic alignment that only uses

26  raw MS2 data from targeted proteomics or DIA experiments for retention time alignment. To

11

1    optimize the performance of our algorithm, we have investigated the effect of algorithmic

2    parameters on the accuracy of the alignment of runs from validation dataset[5]. First, we evaluate

3    the performance for different similarity measures of chromatogram groups. The dot product

4    masked with spectral angle (see the Supplemental Section S5) as a similarity measure provides

5    the highest fraction of peptides aligned for all 120 possible run-pairs (Fig. 2$a$). Within an RT

6    error tolerance of half peak-width (here: 15.3 sec), this similarity measure aligns 94.33% of

7    annotated peaks with the highest area under the curve of all approaches (see the supplemental

8    Table 7 and 8).

9          We then investigated the effect of gap penalty used in dynamic programming. In

10   DIAlignR, the gap penalty was calculated heuristically as a fixed quantile value of the

11   distribution of similarity scores. We find that the selection of quantile value does not have a

12   considerable impact on the percentage of peaks aligned within certain RT tolerance (Fig. 2$b$).

13   From the figure, 20th to 90th quantile values yield approximate 95.6% of aligned peaks within half

14   peak-width. The effect of gapQuantile is less pronounced for wider RT tolerance. For further

15   analysis, the 65th quantile is selected as the base gap penalty for chromatogram alignment. For

16   the affine gap penalty, a gap opening factor of 0.125 was used, while a gap extension factor of

17   40 was considered (see the Supplemental Section S5).

18          We next investigated the impact of the number of transitions on the alignment(see the

19   Supplemental Section S5 and supplemental Figure 10$a$). We find that as number of fragment-

20   ions increases, alignment accuracy improves with best alignment 94.3% from using all library

21   fragment-ions compared to 88.7% with only one transition (Supplemental Table 15). However,

22   the largest improvement was going from 1 transition (88.7%) to two transitions (92.8%) with

23   smaller improvements for using all transitions. This indicates single transition is not sufficient for

24   this algorithm, but the relative gain is more modest after two transitions.

25          Our algorithm can constrain the similarity matrix using a global alignment function.

26   Constraining the alignment in a certain window (given by RSEdistFactor) around the global fit

1 derived from "anchor peptides" improves the alignment accuracy. We have observed that with a

2 constrained similarity matrix 95.4% peaks get aligned compared to 94.3% with non-constrained

3 one (see the supplemental Figure 10*b*). An example of such alignment is shown in Fig. 2*c* and

4 2*d*, in which the similarity matrix has two high similarity hot-spots. By constraining the alignment

5 inside the dashed region in Fig. 2*d*, the alignment path goes through the correct hot spot. With

6 the unconstrained similarity matrix, an incorrect alignment results as shown in the supplemental

7 Figure 11.

8

9 *Validation using "gold standard" reference dataset*

10 Using validation dataset, we have compared DIAlignR to current alignment methods. In

11 terms of number of peptides aligned and alignment precision, chromatogram alignment

12 outperforms LOESS and linear regression methods (see the Fig. 3*a* and Table I). On the *S.*

13 *pyogenes* benchmark dataset, DIAlignR decreases error rates by 1.8-fold compared to the

14 state-of-the-art methods. Cumulatively, chromatographic alignment only mis-aligns 4.3% of all

15 peaks within 15.3 seconds (half peak width) of the true RT compared to 7.9% for LOESS (while

16 LOESS with default parameters mis-aligns 22.8% of all peaks and linear regression mis-aligns

17 44.8%; see Fig. 3*a*).

18 We next investigated the effect of experimental perturbation on the performance of the

19 alignment method. We compare within-condition alignments with between-condition alignments

20 (in the validation dataset, the conditions are 0% and 10% human plasma added to *S. pyogenes*

21 during growth). For both alignment methods, we observed decreased performance for between-

22 condition alignments compared to within-condition alignments (Fig. 3*b*). However, the

23 performance drop of the LOESS method (4.93%) is substantially larger than the corresponding

24 performance drop of DIAlignR (2.7%), indicating increased robustness to sample heterogeneity

25 for DIAlignR.

1      To evaluate the consistency of alignment approaches across multiple run-pairs, we have

2      computed the number of correct aligned peaks (defined as instances where the alignment is

3      correct within half peak-width) for each run-pair. This distribution is shifted towards the right with

4      low standard deviation for chromatogram alignment method compared to LOESS, indicating that

5      the former is consistent in its performance (Fig. 3*c*). In terms of the precision of the alignment,

6      chromatogram alignment consistently performs better than global alignment methods as the

7      former has higher area under the cumulative peptide frequency curve for each run-pair (higher

8      AUC for 120 out of 120 pairs, see the supplemental Figure 12*c*). Similarly, we observe a larger

9      RT variation (standard deviation = 18.45 sec) with the LOESS approach which chromatogram

10      alignment corrects satisfactorily with the standard deviation being 11.68 sec (Fig. 3*d* and

11      supplemental Fig. 12*a*, *b*). We conclude that on the validation dataset, DIAlignR performs

12      consistently better in terms of accuracy of alignment and number of correctly aligned peaks

13      across a range of different RT cutoffs and LC-MS/MS runs.

14      Next, we were interested in how the global differences between the two methods

15      translate to individual alignments. We, therefore, computed the alignment error for each

16      pairwise alignment of each peptide (total 49,505 alignments) and found that chromatographic

17      alignment outperforms LOESS in 4.7% of all cases, whereas, LOESS achieves better results in

18      1.1% cases, while comparable performance was achieved in the remaining 94.2% cases (see

19      the Supplemental Fig. 13*b*). On average, DIAlignR has reduced the RT error by 2.3 seconds

20      with a median of 1.7 seconds (see the supplemental Fig. 13*c*). Overall, our method aligns 47.3k

21      peaks correctly compared to 45.6k by an optimized LOESS within half peak-width (15.3 sec).

22      However, in general we observed that on the validation dataset both methods perform with

23      similar consistency which may be due to the low complexity of a bacterial sample and the high

24      homogeneity of the data as it was acquired within two consecutive days on the same LC

25      column.

26

1    *Application to large-scale heterogeneous human plasma measurements*

2    After demonstrating consistently improved performance on the *S. pyogenes* validation

3    dataset, we investigated the performance of our algorithm on a large-scale SWATH-MS

4    experiments on human plasma. These experiments provided a more challenging dataset as the

5    data was acquired over the period of six months with an intermittent repair of the instrument and

6    replacement of the old column (column1) with a new column (column2). Two LC-MS/MS runs

7    from each of the 12 batches were selected at random and 406 peptides were used for testing

8    our algorithm. Since we did not have manually validate peaks, high confidence peak groups (q-

9    value $< 10^{-3}$) were used instead as validation set.

10    Comparing our chromatogram alignment algorithm (DIAlignR) with the LOESS method

11    on a highly heterogeneous human plasma dataset, we found that our approach aligns 97.92%

12    of peaks compared to 76.03% by LOESS with a maximal error of 20 seconds (half

13    chromatographic peak-width) as depicted in Fig. 4*a*. All tested 276 pairwise alignments shown

14    improved performance using chromatographic alignment (see Supplemental Fig. 15). Next, we

15    were interested in the performance of our method on the alignment of runs acquired on the two

16    different columns. We find that for runs acquired on different columns, chromatogram alignment

17    method aligns 97.7% (compared to 97.84% within-column alignment) of peaks compared to

18    63.38% (89.06% for within-column alignment) by LOESS method (Fig. 4*b*), suggesting that

19    DIAlignR retains performance even for highly heterogeneous datasets while the LOESS

20    approach does not. Specifically, we find not only that DIAlignR outperforms LOESS on between-

21    column alignments, but that the performance loss for DIAlignR is much less pronounced than for

22    LOESS compared to within-column alignments. After validating the performance of

23    chromatogram alignment cumulatively, we decided to investigate its consistency across

24    individual run-pair alignments. Fig. 4*c* presents the distribution of the number of peaks aligned in

25    all 276 pairs. DIAlignR is capable of aligning 400 peaks on average within half-peak width (while

1    LOESS aligns 309 peaks on average), a 29% improvement. This indicates substantial alignment

2    error when using LOESS, which could be reduced drastically by DIAlignR.

3        To validate the performance on individual alignment, we further computed the alignment

4    error for each pairwise alignment for every peptide. The standard deviation of alignment error

5    for LOESS was 22.91 sec compared to 13.7 sec for DIAlignR (Fig. 4*d*). This indicates the higher

6    precision of RT alignment with our approach. Out of 112,056 alignments, we found that

7    DIAlignR outperforms LOESS in 23% of all cases and performs similarly in 76% cases (see the

8    Supplemental Figure 14 - while performing worse in 1.13% cases).  Upon manual validation,

9    several of these worse-performing peaks were found to be due to wrong annotation by

10   OpenSWATH (Supplemental Figure 18). Thus, testing of chromatogram alignment approach on

11   the heterogeneous human plasma dataset again validates its consistent and improved RT

12   alignment performance.

13

14   *Switching of peptide elution order*

15       In liquid chromatography, retention time drift is often observed from one run to another

16   run. However, we were interested whether this drift may be variable for different peptides and

17   thus could result in reversal of retention order[15]. In such a scenario, two peptides which are

18   eluting in order in one run may reverse their elution order in another run. Since our approach

19   does not make an assumption of order preservation of peptide elution and facilitate independent

20   alignment, we hypothesized that DIAlignR would be capable of uncovering instances of non-

21   order preserving chromatographic alignment. Specifically, we analyzed peptide pairs that switch

22   elution order from the heterogeneous and distant plasma runs.

23       To confirm the alignment for such peak-switching cases by chromatogram alignment

24   algorithm, we have specifically looked at the alignment of the pair "run4_run23" as it has the

25   highest number of peak switching pairs. run4 is part of batch V4, acquired on February 28[th],

26   2017 whereas run23 is from batch M3, acquired on July 20[th], 2017. The LOESS fitting from

16

1    common high scoring training peptides for this pair is presented in Fig. 5*a*. Most of the test

2    peptides are scattered around the global fit line, instead of being directly on the line. This graph

3    quickly suggests 407 peptide pairs (one from either side of the line) comprising of 237 peptides

4    that have switched their elution order out of 406 peptides (see the supplemental Section S6).

5    We thus found that overall, 58.4% of peptides were involved in at least one event of non-order

6    preserving elution.

7        One of the peak switching cases is presented in Fig 5*b*. In run4 peptide AQLVDMK/2

8    elutes after HYDGSYSTFGER/2, whereas in run23 the elution order has been reversed. Both

9    peptides have seen positive RT drift in run23 from run4, however, HYDGSYSTFGER/2 shifts of

10   1070-850 = 270 seconds whereas AQLVDMK/2 drifts of only 1050-900 = 150 seconds. This

11   varying RT drift between two runs has caused the peptides to elute in different order. The

12   peptide-pair cannot be aligned with a global alignment approach, which in the best-case

13   scenario will be off by 120 seconds -- however, our chromatogram alignment method has

14   mapped the peaks correctly from run4 to run23 (see the supplemental Figure 17).

15       We, further, calculate the cumulative fraction of peptides aligned for pair "run4_run23"

16   (Fig. 5*c*). Chromatogram alignment correctly aligns 98% peaks compared to LOESS which is

17   able to align only 37.93%, thus DIAlignR is capable of decreasing the error by up to 30-fold.

18   Eight peaks, which are not correctly aligned by chromatogram alignment, are further inspected

19   visually by the authors and are found to be the cases of mis-annotation by OpenSWATH, mainly

20   due to the post-translational modifications (see the supplemental Section S7 and supplemental

21   Figure 18).

22

23   **Discussion:**

24       Correcting for retention time drift and aligning retention times between LC-MS/MS runs

25   has been a long-standing problem in proteomics and it has become of particular importance as

26   proteomics moves towards large-scale analysis of human cohorts. However, most efforts so far

17

1    have focussed on MS1 data and few algorithms are available that can exploit the full information

2    present in MS2 spectra produced by targeted methods or DIA / SWATH-MS.

3        In this paper, we have presented a novel algorithm that uses the raw fragment-ion

4    chromatograms directly to perform retention time alignment for targeted proteomics and DIA

5    data. Our algorithm uses XICs to map peaks across a pair of runs and improves accuracy

6    compared to current state-of-the-art methods. We have, furthermore, extended the algorithm

7    and implemented a hybrid approach, which uses a feature-based global alignment to condition

8    the similarity matrix $s$ that led to further gains in accuracy (see the Supplemental Fig. 10$b$). This

9    hybrid approach provides the best of both worlds with a flexible "knob" which allows the user to

10   either put more focus on global features or rely more on local information. To our knowledge,

11   researchers have not yet explored dynamic programming-based alignment on raw fragment-ion-

12   chromatograms. The dynamic programming approach is essential for obtaining a non-linear (or

13   gapped) alignment as distant runs also have varying drift even for local peaks. Using LOESS to

14   partially constrain the alignment makes our algorithm more stable and provides the robustness

15   of global alignment methods.

16       We have shown that on a "gold-standard" validation dataset, DIAlignR consistently

17   outperforms a global alignment method (using either linear or non-linear approaches), the

18   current state-of-the-art (Supplemental Figure 12$c$). We have observed that alignment accuracy

19   increases almost 4% if a precursor has two transitions instead of one (Supplemental Figure

20   10$a$), We get additional 1.5% aligned precursors by using all six fragment-ions which also

21   corresponds to recommended guidelines[4]. Since, the DIAlignR approach puts more emphasis

22   on local data, we also observe instances of over-fitting where the global alignment function

23   produces better alignment (Supplemental Figure 13$b$). Overall, however, we see increased

24   performance and the DIAlignR algorithm can decrease error rates from 7.9% to 4.3%.

25   Interestingly, we find that our method is less sensitive to changes in chromatographic condition

26   or sample matrix than global alignment approaches (Fig. 3$b$).

1    This finding led us to speculate that the novel chromatographic alignment would be less

2    sensitive to heterogeneity in sample composition and chromatographic condition in large-scale

3    studies. We have tested our algorithm on a large-scale SWATH-MS experiment of human

4    plasma acquired over several months. On this heterogeneous dataset, DIAlignR reduces RT

5    alignment error from 24% to 2%, which is a significant improvement over current state-of-the-art

6    methods. Our approach has outperformed other methods and has consistently mapped the

7    highest number of peaks within half-peak width irrespective of acquisition time interval, column

8    change or instrument repair between two runs (Fig. 4*b*). DIAlignR improves retention time

9    alignment accuracy which has the potential to improve peak-group identification and

10   quantification through downstream tools. We have manually identified an example of wrongly

11   aligned peak-group as shown in Supplemental Figure 19. We have also observed that in the

12   case of a peak being outside of an extracted chromatogram, our method is able to map

13   retention time outside of it as our hybrid approach also uses global alignment (see

14   Supplemental Figure 18). Chromatograms can then be re-extracted and be used to correctly

15   annotate peaks. Thus, this method can further be employed to extract chromatograms by

16   OpenSWATH and other tools.

17   The improvement in retention time alignment comes with non-negligible computation

18   cost. For an alignment of 10-min wide XICs with 3.4 sec cycle time, DIAlignR takes on average

19   0.16 seconds. Thus, pairwise alignment of selected 437 peptides takes approx 1 minute per

20   run-pair. The highest cost in the alignment of each peptide is the calculation of the alignment

21   path using dynamic programming. However, this problem scales linearly with the number of

22   peptides in the library and can be easily parallelized on a computing cluster.

23   We believe that our approach is most useful for large-scale heterogeneous targeted

24   proteomics studies where runs are acquired by different personnel and data is collected over

25   several months or even years. Applying a single mapping function in such experiments

26   becomes a very challenging task considering the switching of elution order of peptides. Global

1    alignment functions, being monotonic in nature, assume chronological order of peptide elution

2    and, therefore, cannot align switched peptides. Hence, we have observed substantial

3    underfitting with the global function and an overall reduction of error using DIAlignR. Our hybrid

4    approach aligns these peptides accurately as it mostly relies on additional dimensions of

5    fragment-ion m/z to align peaks. It is possible that switching peptides may share many

6    fragment-ions, however, this is very rare scenario if library is designed carefully[4] and in such

7    cases our method will perform no worse than global alignment methods.

8         Accurate RT alignment has multiple uses in the application of mass spectrometry-based

9    proteomics for large-scale systems biology studies. Correct identification and improved

10   quantitation of large number of analytes are few of them. This seems intuitive as most

11   quantitative approaches currently available at least to some degree rely on accurate retention

12   time alignment. We present a tool that can align retention times of DIA data by establishing

13   correspondence between analytes across large number of samples, making DIA amenable for

14   multi-center and longitudinal studies. We also expect that this tool can be utilized by existing

15   proteomics software to streamline analyte identification and improve the quantification.

16

17   **Acknowledgements:**

21

22   **Author contributions:**

23   S.G. designed and wrote code, performed data analysis and produced the figures. S.A. and

24   W.Z. performed human plasma related experiments, acquired MS data. H.R. designed and

25   supervised the study. S.G. and H.R. contributed to writing the manuscript. All authors have

26   contributed to the final manuscript.

1

2 **Competing Financial Interest:**

3 The authors declare no competing financial interests.

4

5 **Data Availability:**

6 Raw chromatograms and features extracted by OpenSWATH are available on PeptideAtlas

7 under accession code PASS01280.

8

9 **References**:

10 1. Uzozie, A. C. & Aebersold, R. Advancing translational research and precision medicine with

11 targeted proteomics. *J. Proteomics* (2018). doi:10.1016/j.jprot.2018.02.021

12 2. Surinova, S. *et al.* On the development of plasma protein biomarkers. *J. Proteome Res.* **10,**

13 5–16 (2011).

14 3. Nigjeh, E. N. *et al.* Quantitative Proteomics Based on Optimized Data-Independent

15 Acquisition in Plasma Analysis. *J. Proteome Res.* **16,** 665–676 (2017).

16 4. Schubert, O. T., Röst, H. L., Collins, B. C., Rosenberger, G. & Aebersold, R. Quantitative

17 proteomics: challenges and opportunities in basic and applied research. *Nat. Protoc.* **12,**

18 1289–1294 (2017).

19 5. Röst, H. L. *et al.* TRIC: an automated alignment strategy for reproducible protein

20 quantification in targeted proteomics. *Nat. Methods* **13,** 777–783 (2016).

21 6. Gillet, L. C. *et al.* Targeted data extraction of the MS/MS spectra generated by data-

22 independent acquisition: a new concept for consistent and accurate proteome analysis.

23 *Mol. Cell. Proteomics* **11,** O111.016717 (2012).

24 7. Guo, T. *et al.* Rapid mass spectrometric conversion of tissue biopsy samples into

25 permanent quantitative digital proteome maps. *Nat. Med.* **21,** 407–413 (2015).

8.   Röst, H. L. *et al.* OpenSWATH enables automated, targeted analysis of data-independent acquisition MS data. *Nat. Biotechnol.* **32,** 219–223 (2014).

9.   Tsou, C.-C. *et al.* DIA-Umpire: comprehensive computational framework for data-independent acquisition proteomics. *Nat. Methods* **12,** 258 (2015).

10.  Navarro, P. *et al.* A multicenter study benchmarks software tools for label-free proteome quantification. *Nat. Biotechnol.* **34,** 1130–1136 (2016).

11.  Liu, Y. *et al.* Quantitative variability of 342 plasma proteins in a human twin population. *Mol. Syst. Biol.* **11,** 786 (2015).

12.  Escher, C. *et al.* Using iRT, a normalized retention time for more targeted measurement of peptides. *Proteomics* **12,** 1111–1121 (2012).

13.  Bruderer, R., Bernhardt, O. M., Gandhi, T. & Reiter, L. High-precision iRT prediction in the targeted analysis of data-independent acquisition and its impact on identification and quantitation. *Proteomics* **16,** 2246–2256 (2016).

14.  Smith, R., Ventura, D. & Prince, J. T. LC-MS alignment in theory and practice: a comprehensive algorithmic review. *Brief. Bioinform.* **16,** 104–117 (2015).

15.  Spicer, V., Grigoryan, M., Gotfrid, A., Standing, K. G. & Krokhin, O. V. Predicting retention time shifts associated with variation of the gradient slope in peptide RP-HPLC. *Anal. Chem.* **82,** 9678–9685 (2010).

16.  Wu, L., Amon, S. & Lam, H. A hybrid retention time alignment algorithm for SWATH-MS data. *Proteomics* **16,** 2272–2283 (2016).

17.  Nielsen, N.-P. V., Carstensen, J. M. & Smedsgaard, J. Aligning of single and multiple wavelength chromatographic profiles for chemometric data analysis using correlation optimised warping. *J. Chromatogr. A* **805,** 17–35 (1998).

18.  Bylund, D., Danielsson, R., Malmquist, G. & Markides, K. E. Chromatographic alignment by warping and dynamic programming as a pre-processing tool for PARAFAC modelling of liquid chromatography–mass spectrometry data. *J. Chromatogr. A* **961,** 237–244 (2002).

19. Listgarten, J., Neal, R. M., Roweis, S. T. & Emili, A. Multiple Alignment of Continuous Time Series. in *Advances in Neural Information Processing Systems 17* (eds. Saul, L. K., Weiss, Y. & Bottou, L.) 817–824 (MIT Press, 2005).

20. Sadygov, R. G., Maroto, F. M. & Hühmer, A. F. R. ChromAlign: A two-step algorithmic procedure for time alignment of three-dimensional LC-MS chromatographic surfaces. *Anal. Chem.* **78,** 8207–8217 (2006).

21. Prakash, A. *et al.* Signal maps for mass spectrometry-based comparative proteomics. *Mol. Cell. Proteomics* **5,** 423–432 (2006).

22. Hoffmann, N. & Stoye, J. ChromA: signal-based retention time alignment for chromatography-mass spectrometry data. *Bioinformatics* **25,** 2080–2081 (2009).

23. Christin, C. *et al.* Time alignment algorithms based on selected mass traces for complex LC-MS data. *J. Proteome Res.* **9,** 1483–1495 (2010).

24. Robinson, M. D. *et al.* A dynamic programming approach for the alignment of signal peaks in multiple gas chromatography-mass spectrometry experiments. *BMC Bioinformatics* **8,** 419 (2007).

25. Baran, R. *et al.* MathDAMP: a package for differential analysis of metabolite profiles. *BMC Bioinformatics* **7,** 530 (2006).

26. Wang, J. & Lam, H. Graph-based peak alignment algorithms for multiple liquid chromatography-mass spectrometry datasets. *Bioinformatics* **29,** 2469–2476 (2013).

27. Sandin, M. *et al.* An adaptive alignment algorithm for quality-controlled label-free LC-MS. *Mol. Cell. Proteomics* **12,** 1407–1420 (2013).

28. Prince, J. T. & Marcotte, E. M. Chromatographic alignment of ESI-LC-MS proteomics data sets by ordered bijective interpolated warping. *Anal. Chem.* **78,** 6140–6152 (2006).

29. Hoffmann, N. *et al.* Combining peak- and chromatogram-based retention time alignment algorithms for multiple chromatography-mass spectrometry datasets. *BMC Bioinformatics* **13,** 214 (2012).

1    30.  Chambers, J. M., Hastie, T. J. & Others. *Statistical models in S.* **251,** (Wadsworth &

2         Brooks/Cole Advanced Books & Software Pacific Grove, CA, 1992).

3    31.  Searle, B. C. *et al.* Comprehensive peptide quantification for data independent acquisition

4         mass spectrometry using chromatogram libraries. *bioRxiv* 277822 (2018).

5         doi:10.1101/277822

6    32.  Rosenberger, G. *et al.* Statistical control of peptide and protein error rates in large-scale

7         targeted data-independent acquisition analyses. *Nat. Methods* **14,** 921–927 (2017).

8    33.  Chambers, M. C. *et al.* A cross-platform toolkit for mass spectrometry and proteomics. *Nat.*

9         *Biotechnol.* **30,** 918–920 (2012).

10   34.  Durbin, R., Eddy, S. R., Krogh, A. & Mitchison, G. *Biological sequence analysis:*

11        *probabilistic models of proteins and nucleic acids.* (Cambridge university press, 1998).

1      Table I
2      Summary description of validation and human plasma datasets
3

| | Validation dataset | Human plasma dataset |
|---|---|---|
| Biological Sample | *Streptococcus pyogenes strain SF370* | Plasma from blood samples |
| Mass-spectrometer | SciEX 5600 TripleTOF | SciEX 6600 TripleTOF |
| LC - gradient | Linear | Linear |
| Total run time | 135 minutes | 55 minutes |
| LC-column replaced | No | Yes |
| Mass-spectrometer repaired | No | Yes (Replaced quadrupole after 7 batches) |
| Data acquisition date | 08 August 2012 - 09 August 2012 | 17 February 2017- 20 July 2017 |
| Number of runs acquired | 16 | 975 |
| Number of batches | 1 | 12 |
| Runs selected for alignment | 16 | 24 |
| Total number of run-pairs | 120 | 276 |
| Software used for feature detection and XIC extraction | OpenSWATH | OpenSWATH |
| Number of common precursors selected per run for alignment | 437 | 406 |
| Total number of alignments | 49505 | 112,056 |
| Manual Annotation | Yes (Skyline) | No |

4
5

25

1                                    Table II

2    Below is presented average number of peptides aligned from manually validated *S. Pyogenes*

3    dataset and from *heterogenous human plasma measurements*. For plasma data, peaks with

4    OpenSWATH m-score < 0.001 are used for evaluation of the algorithm.

| Dataset | Method | Average number of peaks aligned within half peak-width | Average number of peaks aligned within one peak-width | Average number of peaks aligned within two peak-width | Total number of peaks per run | Peaks aligned within half peak-width (%) |
|---|---|---|---|---|---|---|
| *Validation* dataset | Chromatogram Alignment | 394.81667 | 407.3 | 410.69167 | 437 | 95.68765 |
| | LOESS | 380.07500 | 406.18333 | 410.81667 | 437 | 92.10249 |
| *Heterogenous human plasma* dataset | Chromatogram Alignment | 397.56159 | 400.50362 | 403.88043 | 406 | 97.92157 |
| | LOESS | 308.67754 | 382.17754 | 403.15217 | 406 | 76.02895 |

5

26

**a**

Intensity

Retention Time (Chrom A)  $t_I$

Intensity

Retention Time (Chrom B)  $t_J$

**b**

$\vec{b}$  $\vec{a}$

$\theta$

$$\vec{a} \cdot \vec{b} = \|a\| \ \|b\| \cos(\theta) = a_1b_1 + a_2b_2 + a_3b_3$$

**c**

$j$

$i$

Similarity score matrix

High

Low

$(I, J)$

**f**

Intensity

Retention Time (Chrom A)  $t_{I'}$

Intensity

Retention Time (Chrom B)  $t_{J'}$

**e**

Matrix M
Matrix A
Matrix B

High

Low

**d**

$j$

$i$

High

Low

$(I, J)$

FIG. 1. **Alignment algorithm for targeted proteomics MS2 chromatograms.** *a*, Fragment-ions chromatograms of a peptide for two runs; *run A* at top and *run B* at bottom. Correct peak, typically, has all library fragment-ions ($n = 3$) coeluting. *b*, similarity between chromatograms of both runs is calculated by dot-product of intensity vector; defined in *n* dimensional space. *c*, outer dot-product of chromatograms provides an *I* x *J* similarity score matrix (*S*). *d*, feature-based complete run alignment is used as an approximate path for alignment. Time points farther from an allowed window in similarity score matrix are penalized by adding negative score. *e*, Affine gap penalty based overlap alignment strategy is employed for calculating best scoring path through the similarity matrix. This dynamic programming based strategy utilizes three matrices for recursively calculating multiple gap length scores. Calculated alignment path is indicated using black arrow. *f*, Chromatograms recreated by mapping intensity back to aligned time path.

**a** Comparison for different similarity measures

Cumulative fraction of peptides vs Retention time difference (in sec)

- Euclidean distance
- Pearson correlation
- 2*spectral angle
- Covariance
- Dot product
- Dot product masked

**b** Effect of gapQuantile value on the number of aligned peaks

Percentage of peaks aligned vs gapQuantile

- Alignment tolerance = Half peak–width
- Alignment tolerance = One peak–width
- Alignment tolerance = Two peak–width

**c** Alignment path through the similarity matrix
for 7481_DGSVSVADSGR/2

High similarity
Vectors A

High similarity
Vectors B

run12 index vs run11 index

**d** LOESS fit between run11 and run12 (span = 0.03 )
run11: Strep0Repl1_R02_SW
run12: Strep10Repl1_R02_SW

run12 time vs run11 time

**F**IG. 2. **Comparison of different similarity measurements, technical parameters and effect of penalizing similarity using global prior on the accuracy of alignment in *S. Pyogenes* dataset.** *a*, Performance of various similarity measures as the cumulative fraction of peptides having error less than RT difference is plotted. *b*, the effect of gap penalty selection using gapQuantile on the percentage of peaks aligned within certain RT difference tolerance is depicted. *c*, The penalized similarity matrix for peptide DGSVSVADSGR/2 between run11 and run12 is presented. From available two high similarity vectors, alignment path passes through high similarity vectors B. *d*, The end-points of extracted ion chromatograms (XICs) for the peptide are shown as green dots. Penalizing similarity gives preference to alignment within a certain window around LOESS fit, depicted as dashed green lines. Here, alignment of high similarity vectors B (solid red circle ●) is preferred over high similarity vectors A (red cross ☐).

**a** Between all 120 pairs

**b** Change in fit with pair-run type

**c** Histogram of number of aligned peaks in 120 pairs
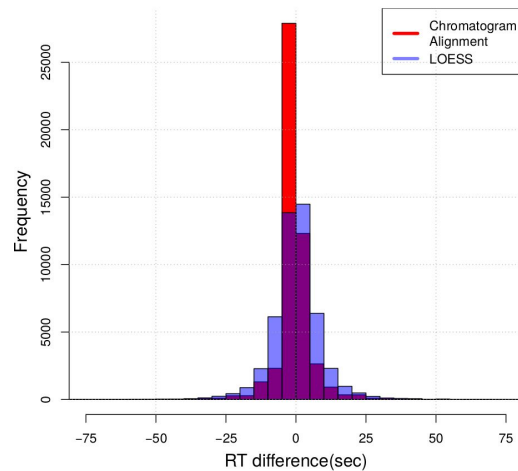
**d** Histogram of RT difference

FIG. 3. **Alignment accuracy of MS2 chromatogram alignment on a validation dataset of 16 runs with manually annotated 437 peak groups in each run.** *a*, cumulative fraction of peptides having error less than RT difference is plotted for all possible C(16,2) = 120 pairs for chromatogram alignment, linear fit and k-nearest neighbor smoothing (LOESS) with and without optimum span. *b*, cumulative fraction of peptides with alignment accuracy is plotted for chromatogram alignment and LOESS for pairs with different biological conditions. Strep0 pair constitutes both 0% plasma runs, Strep10 pair is composed of both 10% plasma runs and Strep0_Strep10 pair have one run with 0% plasma and other with 10% plasma. There are 28 Strep0 pairs, 28 Strep10 pairs and 64 pairs for Strep0_Strep10 case in the validation dataset. *c*, histogram of number of peptides matched within half peak-width for LOESS and chromatogram alignment. *d*, Histogram of retention time (RT) prediction error is plotted for chromatogram alignment and LOESS. RT difference standard deviation for both approaches is 9.56 sec and 10.98 sec, respectively.
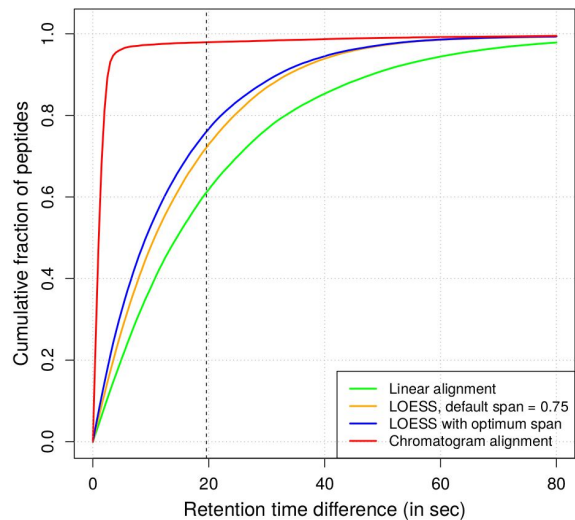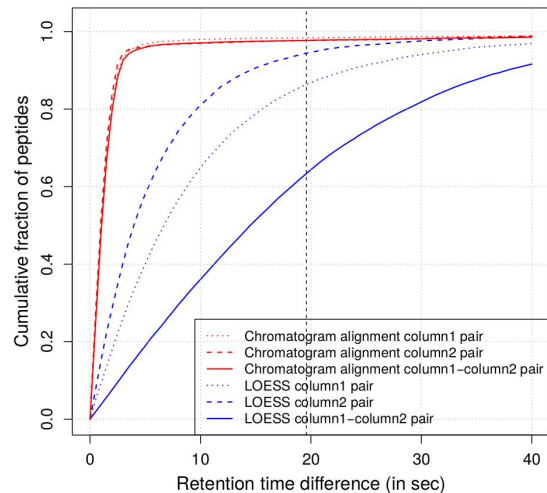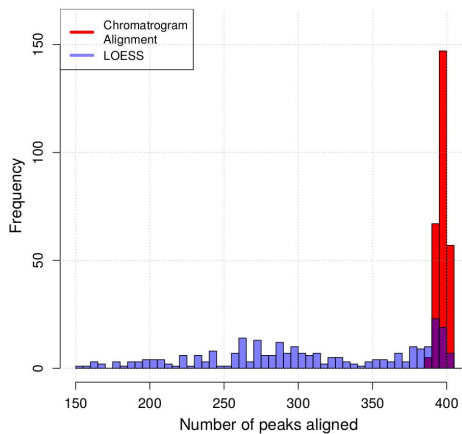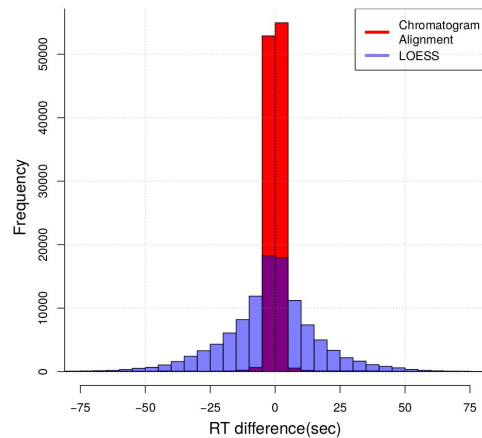
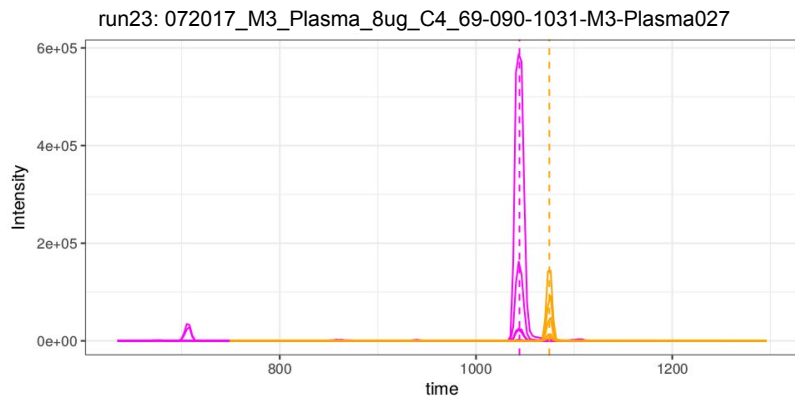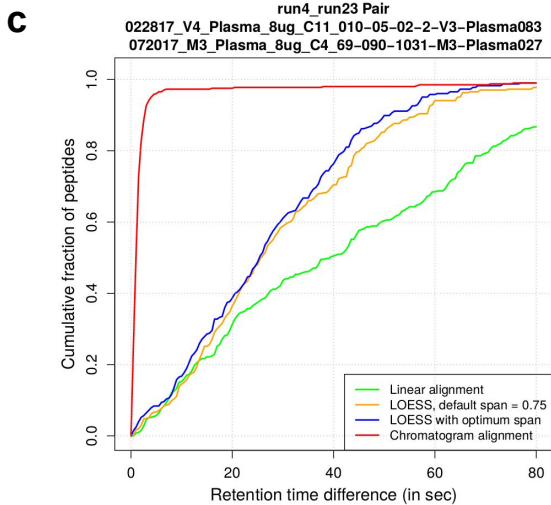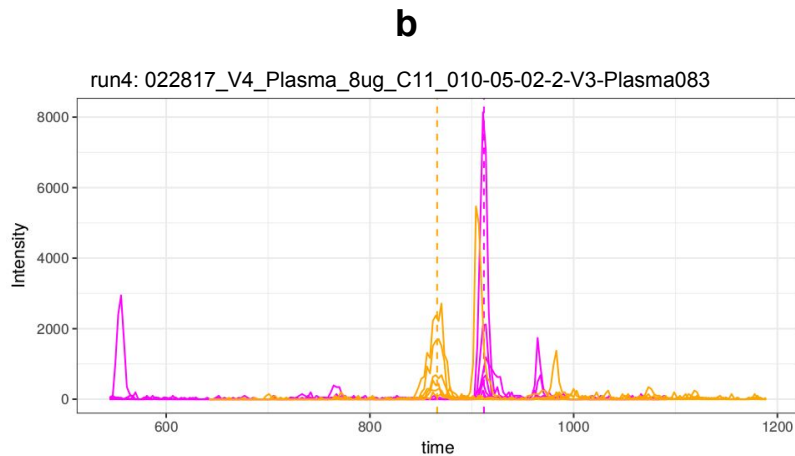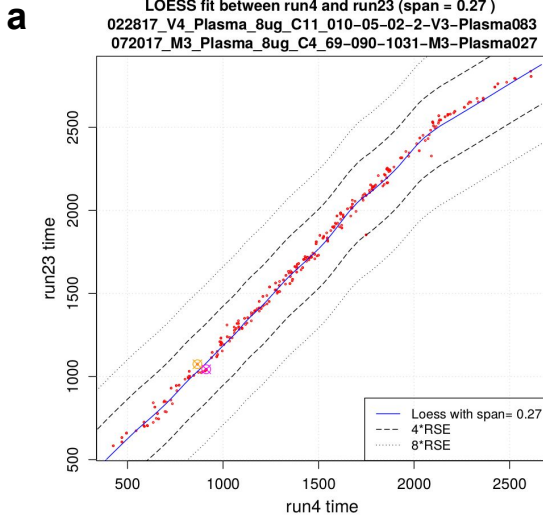**a** Between all 276 pairs

**b** Alignment performance with column change

**c** Histogram of number of aligned peaks in 276 pairs

**d** Histogram of RT difference

FIG. 4. **Alignment accuracy of MS2 chromatogram alignment on 24 runs of clinical plasma measurement dataset annotated with OpenSWATH. 406 peak groups are selected in each run with m-score < 0.001.** *a*, cumulative fraction of peptides having error less than RT difference is plotted for all possible C(24,2) = 276 pairs for chromatogram alignment, linear fit and k-nearest neighbor smoothing (LOESS) with and without optimum span. *b*, cumulative fraction of peptides with alignment accuracy is plotted for chromatogram alignment and LOESS for pairs with different data acquisition conditions. LC column was changed together with quadrupole replacement. 14 runs were acquired on column1 which makes 91 pairs, labeled as "column1". 10 runs were acquired after quadrupole replacement on column2 which results into 45 pairs, labelled as "column2". There are 140 pairs composed of "column1" and "column2" labelled runs; these pairs are labelled as "column1-column2". *c*, histogram of number of peptides matched within half peak-width for LOESS and chromatogram alignment. *d*, Histogram of retention time (RT) prediction error is plotted for chromatogram alignment and LOESS. RT difference standard deviation for both approaches is 22.91 sec and 13.7 sec, respectively.

**a** LOESS fit between run4 and run23 (span = 0.27 )
022817_V4_Plasma_8ug_C11_010-05-02-2-V3-Plasma083
072017_M3_Plasma_8ug_C4_69-090-1031-M3-Plasma027

Loess with span= 0.27
4*RSE
8*RSE

**b**

run4: 022817_V4_Plasma_8ug_C11_010-05-02-2-V3-Plasma083

run23: 072017_M3_Plasma_8ug_C4_69-090-1031-M3-Plasma027

119474_AQLVDMK/2
38644_HYDGSYSTFGER/2

**c** run4_run23 Pair
022817_V4_Plasma_8ug_C11_010-05-02-2-V3-Plasma083
072017_M3_Plasma_8ug_C4_69-090-1031-M3-Plasma027

Linear alignment
LOESS, default span = 0.75
LOESS with optimum span
Chromatogram alignment

**F**IG. 5. **Alignment of 406 peptides in pair *run4 and run23* from clinical plasma measurement dataset.** run4 "022817_V4_Plasma_8ug_C11_010−05−02−2−V3−Plasma083" was acquired on February 28<sup>th</sup>, 2017 whereas run23 "072017_M3_Plasma_8ug_C4_69−090−1031−M3−Plasma027" was acquired on July 20<sup>th</sup>, 2017**.** *a*, LOESS fit between two runs is obtained using confident peaks. Test peptides are shown in red color around the fit line. Span value = 0.27 for fit is obtained by ⅓ cross-validation. Precursors AQLVDMK/2  and HYDGSYSTFGER/2 are shown in magenta and orange circle-cross symbols, respectively. *b*, Two peptides AQLVDMK/2 and HYDGSYSTFGER/2 have their elution order reversed in these runs. This phenomenon makes alignment of peaks theoretically impossible for global monotonic methods. Chromatogram alignment uses fragment-ions as additional dimensions and hence can align them precisely. *c*, fraction of peptides having error less than RT difference is plotted for pair *run4 and run23* for chromatogram alignment, linear fit, k-nearest neighbor smoothing (LOESS) with and without optimum span and without any alignment.