



HAL
open science

CHICKN: extraction of peptide chromatographic elution profiles from large scale mass spectrometry data by means of Wasserstein compressive hierarchical cluster analysis

Olga Permiakova, Romain Guibert, Alexandra Kraut, Thomas Fortin,
Anne-Marie Hesse, Thomas Burger

► To cite this version:

Olga Permiakova, Romain Guibert, Alexandra Kraut, Thomas Fortin, Anne-Marie Hesse, et al.. CHICKN: extraction of peptide chromatographic elution profiles from large scale mass spectrometry data by means of Wasserstein compressive hierarchical cluster analysis. *BMC Bioinformatics*, 2021, 22 (1), 10.1186/s12859-021-03969-0 . hal-03145601

HAL Id: hal-03145601

<https://hal.science/hal-03145601>

Submitted on 7 Oct 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 CHICKN: Extraction of peptide chromatographic
2 elution profiles from large scale mass
3 spectrometry data by means of Wasserstein
4 compressive hierarchical cluster analysis

5 Olga Permiakova, Romain Guibert, Alexandra Kraut, Thomas
Fortin, Anne-Marie Hesse and Thomas Burger
thomas.burger@cea.fr

6 Univ. Grenoble Alpes, CNRS, CEA, Inserm, BGE U1038, 38000
7 grenoble, France

8 **Abstract**

9 **Background:** The clustering of data produced by liquid chromatog-
10 raphy coupled to mass spectrometry analyses (LC-MS data) has recently
11 gained interest to extract meaningful chemical or biological patterns. How-
12 ever, recent instrumental pipelines deliver data which size, dimensionality
13 and expected number of clusters are too large to be processed by classical
14 machine learning algorithms, so that most of the state-of-the-art relies on
15 single pass linkage-based algorithms.

16 **Results:** We propose a clustering algorithm that solves the powerful
17 but computationally demanding kernel k -means objective function in a
18 scalable way. As a result, it can process LC-MS data in an acceptable
19 time on a multicore machine. To do so, we combine three essential fea-
20 tures: a compressive data representation, Nyström approximation and a
21 hierarchical strategy. In addition, we propose new kernels based on opti-
22 mal transport, which interprets as intuitive similarity measures between
23 chromatographic elution profiles.

24 **Conclusions:** Our method, referred to as CHICKN, is evaluated on
25 proteomics data produced in our lab, as well as on benchmark data com-
26 ing from the literature. From a computational viewpoint, it is particularly
27 efficient on raw LC-MS data. From a data analysis viewpoint, it provides
28 clusters which differ from those resulting from state-of-the-art methods,
29 while achieving similar performances. This highlights the complementar-
30 ity of differently principle algorithms to extract the best from complex
31 LC-MS data.

32 **Keywords:** Large-scale cluster analysis; Liquid chromatography; Mass-
33 spectrometry; Proteomics; Wasserstein kernel; Optimal transport

1 Background

2 Liquid chromatography coupled to mass spectrometry (LC-MS) constitute a
3 technological pipeline that has become ubiquitous in various omics investiga-
4 tions, such as proteomics, lipidomics and metabolomics. Over the past decade,
5 the MS throughput has continuously improved, leading to unprecedented data
6 volume production. To date, processing these gigabytes of low level MS sig-
7 nals has become a challenge on its own, for a trade-off between contradictory
8 objectives is sought: On the one hand, one needs to save memory and computa-
9 tional time with efficient encoding, compression and signal cleaning methods [1].
10 On the other hand, one needs to avoid too important preprocessing that sys-
11 tematically smoothes signals of lower magnitudes, as it is now well-established
12 that interesting biological patterns can be found near the noise level [2]. To
13 face this challenge, a recent and efficient investigation path has been to ap-
14 ply cluster analysis to LC-MS data. Cluster analysis refers to a large family
15 of unsupervised statistical learning and multivariate analysis techniques which
16 share a common goal: Aggregating similar data items into clusters, so that
17 within-cluster similarities are larger than between cluster ones. By doing so,
18 it becomes possible to consider the various clusters independently, and thus to
19 reduce the computational footprint without any quality loss. Moreover, as each
20 cluster contains similar data elements, it facilitates the extraction of repetitive
21 but small biological patterns.

22 State of the art

23 To date and contrarily to the presented work, investigations have mainly fo-
24 cused on clustering LC-MS data across the chromatographic (or elution time)
25 dimension, *i.e.* when the data elements are MS spectra: MS2grouper [3, 4],
26 Pep-Miner [5], PepMerger [6], the MS-Clustering / MS-Cluster / Pride-Cluster
27 / spectra-cluster series [7, 8, 9, 10], Bonanza [11], CAMS-RS [12], MaRaClus-
28 ter [13], N-cluster [14], and msCRUSH [15]. All these approaches propose to
29 improve peptide identification by benefiting from the aforementioned trade-off:
30 By grouping similar fragmentation spectra into a consensus representation, one
31 clearly reduces the data volume. Moreover, peaks corresponding to random
32 noise should not reinforce between spectra, while on the contrary, small but
33 chemically consistent peaks should [16].

34 Clustering across the mass-to-charge ratio (m/z) dimension, *i.e.* when the
35 data elements are chromatographic profiles (depicting the signal changes along
36 the elution time at a given m/z value), is also insightful for many reasons:
37 First, it proposes an original framework to construct and extract precursor
38 ion chromatograms, which integration is essential for quantitative analysis [17].
39 Second, cluster centroids naturally provide consensus elution profiles which are
40 of interest for retention time alignment [18]. Finally, elution profiles are also
41 essential to disentangle chimeric spectra [19]. Notably if the clustering is suf-
42 ficiently accurate, it can be insightful to disentangle multiplexed acquisitions
43 (*e.g.* Data Independent Acquisition [20], or DIA), without relying on spectral

1 libraries [21, 22]. To date, these practical problems have been tackled in the
2 proteomics literature by applying various heuristics which differ to some extent
3 from the cluster analysis framework. For instance, in DIA-Umpire [23], peptide
4 fragments’ elution profiles are clustered according to their correlations with pre-
5 cursor profiles, so that formally, the approach is more that of classification (*i.e.*
6 supervised) than of clustering (*i.e.* unsupervised). Similarly, in many quan-
7 tification algorithms (Maxquant [24], OpenMS [25], MsInspect [26], Xnet [17])
8 cluster analysis aims to extract isotopic envelopes, *i.e.* to group the elution
9 profiles of several isotopes of a given molecule, within a closed neighborhood of
10 m/z values. As a consequence, two identical profiles in different m/z regions
11 are not grouped together. Although this behavior (that will be referred to as
12 the *envelope assumption simplification* in the rest of the article) concurs with
13 the objective of isotopic envelope reconstruction, it makes the heuristic strongly
14 attached to one objective; and non applicable to other cluster analysis problems.
15 In contrast, we believe generic clustering algorithms would also be of interest,
16 as different tuning would make them appropriate to deal with different objec-
17 tives: *e.g.* by adding must-link/must-not-link constraints [27] so as to guide
18 the demultiplexing task as in the DIA-Umpire case; or by incorporating an m/z
19 difference in the similarity definition, in the case of isotopic envelope extraction;
20 and so on.

21 Moreover, a refine analysis of the algorithms underlying all these (either
22 spectrum or chromatogram) clustering techniques let appear a strong filiation
23 between them: All rely on agglomerative and linkage-based methods, be it pre-
24 viously published algorithms (HAC [28, 29], DBSCAN [30] or UPGMA [31]) or
25 *ad-hoc* procedures developed in the specific context of LC-MS data clustering
26 (proposed in MS2grouper, Pep-Miner, PepMerger, the MS-Cluster series, Bo-
27 nanza, CAMS-RS, N-cluster and XNet). Despite their unquestionable efficiency,
28 some diversity would help. Cluster analysis is as much an art as a science [32]
29 and there does not exist such thing as the perfect clustering – at least, on real
30 data. Most of the time, data analysts need to rely on a toolbox of various al-
31 gorithms to extract the best of their data [33]. With this respect, MS-based
32 omics would benefit from differently principled and complementary algorithms
33 which have demonstrated their efficiency in data science [34]. For instance,
34 spectral clustering [35, 36, 37] (which should not be confused with the cluster
35 analysis of mass spectra [38]), mean shift algorithm [39, 40], and variants of the
36 k -medoids [41] and k -means [42, 43] are of prime interest.

37 Finally, one observes a difference between algorithms dedicated to spec-
38 trum clustering and those dedicated to chromatogram clustering: While the
39 former ones are mainly implemented in an independent manner, the latter ones
40 are all embedded in computational pipelines (DIA-Umpire [23], Maxquant [24],
41 OpenMS [25], MsInspect [26]). The only exception is Xnet [17], which makes it
42 a unique literature reference for algorithmic and low-level comparisons. In addi-
43 tion, Xnet is the most recently published algorithm, and it displays interesting
44 performances on a benchmark dataset.

45 In a nutshell, Xnet is a Bayesian algorithm which aims to cluster elution
46 profiles into isotopic envelopes. More precisely, it starts from the construction

1 of a network with chromatograms as nodes. Then, the network is decomposed
2 into preliminary clusters. The edges within each cluster are scored by estimating
3 the likelihood of two parameters: the correlation between chromatograms and
4 their m/z separation. Finally, the edge validation is carried out using the scores
5 and a chromatogram apex match verification. This leads to the final isotopic
6 envelope construction.

7 Xnet has many strengths: First, it is a parameter free clustering method
8 – the number of clusters can be inferred during the learning process. Second,
9 the time complexity of the algorithm is linear with respect to the number of
10 chromatograms in the data. However, it also has weaknesses: First, it cannot
11 work on raw data and requires an important preprocessing step, referred to as
12 *ion chromatogram extraction*, which denoises the LC-MS map and aggregates
13 independent measurements into well-formed *traces* (*i.e.* lists of peak intensities
14 corresponding to a same ion, identified in consecutive mass spectra). Concretely,
15 starting from a raw file, it is first necessary to extract non trivial information
16 and to store them into an input CSV file with the following columns: m/z ratios,
17 retention times, intensities and trace labels. In addition to be time consuming,
18 it can arguably be considered that excluding the trace construction from the
19 algorithm amounts to transferring a bottleneck question to another preliminary
20 processing, or to a human annotator. Second, it strongly relies on the envelope
21 assumption simplification, making it impossible to group elution profiles which
22 m/z difference exceeds a predefined threshold. The third weakness is related
23 to the generalization capabilities: As acknowledged in [17], there is not enough
24 data to accurately train the probability model underlying Xnet, making it nec-
25 essary to complement it with a Bayesian prior. This obviously questions the
26 applicability to datasets that significantly differ from the ones that served to
27 tune the prior. Finally, Xnet does not provide a consensus chromatogram for
28 each cluster: Its output is a CSV file that only assigns a cluster index to each
29 line of the input CSV file.

30 Objectives and contributions

31 The objective of this article is twofold: First is to propose a new cluster analysis
32 pipeline adapted to the challenging problem of clustering multiplexed chromato-
33 graphic profiles resulting from data independent acquisitions. The second ob-
34 jective is to build this pipeline around an algorithm which is not agglomerative
35 and linkage-based. Concretely, we focused on k -means objective function, for
36 two reasons: First, until recently, it was considered by the proteomics commu-
37 nity as non-applicable to data as big as LC-MS data [7], while recent theoretical
38 progresses have made this scaling-up possible [44] (this explains the historical
39 predominance of agglomerative linkage-based clustering, less computationally
40 demanding); Second, k -means can be reformulated to fit the reproducing ker-
41 nel Hilbert space theory [45] (leading to the so-called kernel k -means frame-
42 work [46]), which provides new opportunities to define similarity measures that
43 capture the biochemical specificities of LC-MS data (a challenge that has consis-
44 tently been pinpointed as essential over the last fifteen years [3, 5, 6, 11, 12, 13]).

1 The contributions of this article are the following: First, it introduces the use
2 of Wasserstein-1 (W1) distance (a.k.a. earth mover’s distance, a.k.a. optimal
3 transport distance) to account for similarities between elution profiles. Second,
4 it shows that combining Nyström method and random Fourier features leads
5 to a dramatic data compression level that makes the k -means objective function
6 minimizable on raw and high resolution proteomics data with a multi-core ma-
7 chine. Finally, it demonstrates the applicability and interest of the method to
8 process proteomics data from DIA experiments.

9 **Methods**

10 **Materials**

11 To conduct our study, we have relied on three datasets. The first one, hereafter
12 referred to as UPS2GT, is a publicly available dataset [23]. To be used as a
13 benchmark for Xnet, this dataset had been preprocessed and manually anno-
14 tated with isotopic envelopes that can serve as ground truth [47]. Moreover,
15 the data had been converted into *centroid* mode, *i.e.* a compressed version of
16 the original *profile* data. In the *profile* mode, each peak of the mass spectrum
17 is represented by intensities reported for several consecutive m/z values, so as
18 to account for the measurement imprecision. In contrast, the *centroid* mode
19 summarises all the values of the profile mode into a single m/z value, located
20 at the center of the measurement distribution. It leads to significantly smaller
21 memory footprint, at the price of blurring the differences between true signal
22 and noise.

23 The second dataset, hereafter referred to as Ecoli-DIA, is the raw output of
24 a DIA analysis of an *Escherichia Coli* sample (containing over 15,000 peptides¹
25 which signals are multiplexed). To avoid any distortion or information loss, it
26 was stored using the profile mode. The resulting file has an important memory
27 footprint of 3.6 GB. Thus, even though chromatogram clustering operates on
28 fraction of the data only (the so-called MS1 acquisitions, see Ecoli datasets:
29 Data preparation section), it requires adapted software tools and methods.

30 Finally, to account for the rapid increment of data size in proteomics (re-
31 sulting from using ever longer LC and ever more resolved MS acquisitions),
32 we have considered a third dataset, exactly similar to the Ecoli-DIA dataset,
33 but acquired as Full-MS instead of as DIA. This means that 100% of the acqui-
34 sition time was dedicated to MS1 signals, so as to mimick the extraction of a
35 much larger DIA dataset resulting from more time- and m/z -resolved acqui-
36 sitions. This so-called Ecoli-FMS dataset has a memory footprint of 3.2 GB.
37 Even though of equivalent size, this dataset is in fact 16 bigger than Ecoli-DIA
38 (four times more MS1 spectra which are four times more resolved), see Ecoli
39 datasets: Data preparation section.

¹We consider that a peptide is characterized by a triplet: its amino acid sequence, a list of post-translational modifications and their localization on the sequence. Accordingly, different isotope measurements can be grouped into a single peptide definition.

1 **UPS2GT benchmark dataset**

2 The UPS2GT dataset [47] resulted from the liquid chromatography coupled to
3 mass spectrometry analysis of 48 human proteins of the Proteomics Dynamic
4 Range Standard (UPS2) on a AB Sciex TripleTOF 5600 instrument using data
5 dependent acquisition with an MS1 ion accumulation time of 250 ms [23].

6 The 28,568,990 detected points in the resulting LC-MS map were anno-
7 tated according to their intensity value, either as informative or as noisy. Over
8 1,2 million informative points were segmented into 57,140 extracted ion chro-
9 matograms referred to as *traces*. Then, the traces were grouped into 14,076
10 isotopic envelopes. These envelopes constitute the dataset ground truth (there-
11 fore, the objective of the clustering task would be to re-build the envelopes from
12 the traces). The final fully annotated data were stored in a CSV file, where each
13 row depicts one LC-MS point with four pieces of information: its mass to charge
14 ratio, retention time, intensity, trace label and envelope label. The points that
15 were assumed noise were given -1 or 0 as trace label.

16 **Ecoli datasets: wet-lab analysis**

17 *Escherichia Coli* bacteria were lysed with BugBuster reagent (Novagen, final
18 protein concentration $1\mu\text{g}/\mu\text{L}$). Around $560\mu\text{g}$ of proteins were stacked in the
19 top of a 4 - 12% NuPAGE ZOOM gel (Life Technologies) and stained with R-
20 250 Coomassie blue. Gel was manually cut in pieces before being washed by six
21 alternative and successive incubations in 25 mM NH_4HCO_3 for 15 min, followed
22 by 25 mM NH_4HCO_3 containing 50% (v/v) acetonitrile. Gel pieces were then
23 dehydrated with 100% acetonitrile and incubated with 10 mM DTT in 25 mM
24 NH_4HCO_3 for 45 min at 56 °C and with 55 mM iodoacetamide in 25 mM
25 NH_4HCO_3 for 35 min in the dark. Alkylation was stopped by the addition of
26 10 mM DTT in 25 mM NH_4HCO_3 (incubation for 10 min). Gel pieces were then
27 washed again by incubation in 25 mM NH_4HCO_3 followed by dehydration with
28 100% acetonitrile. Modified trypsin (Promega, sequencing grade) in 25 mM
29 NH_4HCO_3 was added to the dehydrated gel pieces for incubation at 37 °C
30 overnight. Peptides were extracted from gel pieces in three sequential extraction
31 steps (each 15 min) in 30 μL of 50% acetonitrile, 30 μL of 5% formic acid, and
32 finally 30 μL of 100% acetonitrile. The pooled supernatants were aliquoted and
33 dried under vacuum.

34 The dried extracted peptides were resuspended in 5% acetonitrile and 0.1%
35 trifluoroacetic acid and 500ng were analyzed by online nanoliquid chromatogra-
36 phy coupled to tandem mass spectrometry (LC-MS/MS) (Ultimate 3000 RSLC-
37 nano and the Q-Exactive HF, Thermo Fisher Scientific). Peptides were sampled
38 on a 300 μm 5mm PepMap C18 precolumn (Thermo Fisher Scientific) and sep-
39 arated on a 75 μm 250 mm C18 column (Reprosil-Pur 120 C18-AQ, 1.9 μm , Dr.
40 Maisch HPLC GmbH). The nano-LC method consisted of a 120 minute multi-
41 linear gradient at a flow rate of 300 nl/min, ranging from 5 to 41% acetonitrile
42 in 0.1% formic acid. The spray voltage was set at 2 kV and the heated capillary
43 was adjusted to 270°C. For the Ecoli-FMS dataset, survey full-scan MS spectra

1 (m/z from 400 to 1,400) were acquired with a resolution of 240,000 after the
 2 accumulation of $3 \cdot 10^6$ ions (maximum filling time 200 ms). For the Ecoli-DIA
 3 dataset, survey full-scan MS spectra (m/z from 400 to 1,400) were acquired with
 4 a resolution of 60,000 after the accumulation of $3 \cdot 10^6$ ions (maximum filling
 5 time 200 ms) and 30 successive DIA scans were acquired with a 33Th width and
 6 a resolution of 30,000 after the accumulation of $2 \cdot 10^5$ ions (maximum filling
 7 time set to auto). The HCD collision energy was set to 30%. MS data were
 8 acquired using the software Xcalibur (Thermo Fisher Scientific).

9 **Ecoli datasets: Data preparation**

10 The output of the LC-MS/MS experiments were converted from the proprietary
 11 RAW format into mzXML files using ProteoWizard [48]. It led to files of 11.4
 12 GB (Ecoli-DIA) and of 10.2 GB (Ecoli-FMS), containing several pieces of infor-
 13 mation: discretized spectra under the form of coupled lists of m/z and intensity
 14 values; as well as metadata about the experiment (number of spectra, retention
 15 time range, etc).

16 In the case of the Ecoli-FMS dataset, all the spectra are peptide mass spec-
 17 tra, also termed MS1. However, the Ecoli-DIA datasets contains two types of
 18 spectra: precursor spectra (MS1) and fragmentation spectra (MS2). Thus, to
 19 work on the elution profiles, we have extracted the MS1 signals from the Ecoli-
 20 DIA file. Then, for both files, we have reconstructed chromatographic signals
 21 from MS1 spectrum intensities. As the proposed method aims to work on data
 22 as raw as possible (*i.e.* without preliminary denoising, smoothing and so on),
 23 we converted each mzXML file into an intensity matrix such as the ones of
 24 Figure 1A (Ecoli-DIA) and of Additional File 1 (Ecoli-FMS), where each row
 25 corresponds to a spectrum and each column to an elution profile (despite pos-
 26 sible m/z fluctuations that may hamper the signal continuity). We concretely
 27 constructed each data matrix using the LC-MS analysis time-stamps and a
 28 non-uniform sampling of the m/z range (see Additional File 2 for a detailed
 29 description). Concretely, the resampled m/z values are given by the following
 30 recursive formula:

$$m_{i+1} - m_i = \frac{0.015}{Res_{EXP}} m_i^{\frac{3}{2}}, \quad (1)$$

31 where m_i is the i^{th} sampled m/z value and Res_{EXP} is the instrument resolution
 32 used in the experiment ($Res_{FMS} = 240,000$ and $Res_{DIA} = 60,000$). Finally, we
 33 have linearly interpolated the intensity values at each node m_i of the grid:

$$I_i = I_{\text{left}} + (m_i - m_{\text{left}}) \cdot \frac{I_{\text{right}} - I_{\text{left}}}{m_{\text{right}} - m_{\text{left}}}, \quad (2)$$

34 where m and I pairs with sub-indexes "left", "right" refer the left and right
 35 neighboring peaks. This is followed by the deletion of the few empty columns.
 36 The resulting Ecoli-DIA data matrix is depicted in Figure 1A: it contains around
 37 3,300 rows and 190,000 columns and it has a footprint of 4.8 GB. As expected,
 38 the Ecoli-FMS data matrix (Additional File 1) is bigger: 14,000 rows, 700,000

1 columns and 82 GB. The bar plots in the margins of both figures represent the
2 intensity distribution across the matrix columns and rows. They show that the
3 Ecoli-FMS and Ecoli-DIA matrices have the same structure and intensity range,
4 despite different size.

5 Methodology overview

6 The proposed methodology is composed of three consecutive parts, hereafter
7 detailed:

8 1. Profile similarity definition:

9 As frequently discussed in the literature [3, 5, 6, 11, 12, 13], the choice
10 of a similarity measure that reflects the biochemical semantics of LC-MS
11 data is essential to achieve efficient processing. In this article, we relied on
12 *Wasserstein-1 distance* [49, 50][51] (or W1, detailed in the Metric choice
13 section) and we transformed it into a similarity by applying a negative
14 exponential function: If x_i and x_j are two chromatograms (or columns
15 from the data matrix), their similarity thus reads:

$$k(x_i, x_j) = e^{-\gamma \cdot [d_{W_1}(x_i, x_j)]^p} \quad (3)$$

16 where d_{W_1} is the W1 distance and where γ is a neighborhood parameter,
17 which tuning authorizes up/down scaling the similarity values. The use
18 of a similarity measure of the form of a negative exponential of a distance
19 is convenient, since it makes it possible to apply the *kernel trick* [52] (see
20 Kernel trick section), *i.e.* to apply a machine learning algorithm as if it
21 were operating in a so-called *feature space* (depicting a non-linear data
22 transform which respects the semantic of the chosen similarity measure).

23 2. **Data compression:** Applying the kernel trick can be rather computa-
24 tionally demanding: For a dataset of size N , it requires the computation
25 of a kernel (or similarity) matrix of size $N \times N$. Thus, with between 10^5
26 and 10^6 chromatograms in the Ecoli datasets, computing and storing the
27 kernel matrix is simply not tractable. The purpose of *Nyström method* [53]
28 (see Nyström approximation section) is to replace the kernel matrix by a
29 low rank approximation, as illustrated in Figure 1B. By relying only on
30 the similarities between each data element and a randomly selected sub-
31 set, it provides a dramatic reduction of the computational burden at the
32 price of a small and controlled loss of accuracy. Even though Nyström
33 approximation allows for an efficient computation of the kernel matrix, it
34 does not accelerate the clustering algorithm itself, which requires multiple
35 traversing of the entire dataset (*i.e.* N elements). To cope for this, it
36 has recently been proposed in the compressive learning framework [54] to
37 summarize the entire dataset by a relatively small vector of fixed size, re-
38 ferred to as *data sketch*, and to have the algorithm operating on his sketch
39 only, irrespective of the original data. Concretely, we built the data sketch
40 as an average of *random Fourier features* of the chromatographic profiles
41 in the feature space (see Random Fourier feature sketching section).

1 **3. Cluster and centroid definitions:** Lloyd algorithm [55] (*i.e.* the most
2 classical algorithm to cluster data according to the k -means objective
3 function) cannot directly be applied on sketched data. Fortunately, it
4 is possible to rely on the *Compressive k -means* (CKM) algorithm pro-
5 posed in [56] (see Cluster computations section). However, CKM only
6 returns a set of cluster centroids and does not cluster the data *per se*.
7 Therefore, traversing the entire (original) dataset to perform the *assign-*
8 *ment* of each chromatogram to its closest centroid (according to the W1
9 distance) is necessary (see Cluster assignment section). CKM complexity
10 does not depend on the original data size (as it operates on the data sketch)
11 which makes it well-scalable. However, its complexity grows rapidly with
12 the number of clusters, which is an issue as thousands of clusters can be
13 sought in LC-MS data. To cope for this, we implemented a *hierarchical*
14 *clustering scheme*, where each cluster is recursively divided into a small
15 number of sub-clusters until the desired number of clusters is obtained (see
16 Cluster assignment section). This procedure provides a set of clusters with
17 centroids only defined in the feature space. To recover the corresponding
18 consensus chromatograms, one has to solve a *pre-image problem*. We prac-
19 tically did so by computing the mean of the elution profiles neighboring
20 each centroid (see Pre-image computation section).

21 To the best of the authors' knowledge, this work is the first one to combine
22 Nyström method and compressive learning with random Fourier features on
23 a problem as difficult as the clustering of LC-MS data, which combines high-
24 dimensionality and a very large number of potential clusters in addition to the
25 traditional difficulties of raw biological data (non-linearities, low signal-to-noise
26 ratio, *etc.*). From this point on, we refer to the proposed method as CHICKN
27 (standing for Chromatogram HIERarchical Compressive K-means with Nyström
28 approximation).

29 **Profile similarity definition**

30 **Metric choice**

31 Originally, the Wasserstein-1 (W1) metric was defined to compute optimal trans-
32 port strategies, which explains why it is also referred to as the *earth mover's*
33 *distance*. It has witnessed a recent gain of interest in machine learning as an ef-
34 ficient way to measure a distance between two probability distributions [57, 58]:
35 Essentially, if one sees probability distributions as earth heaps, the most energy
36 efficient way to move one earth heap in place of the other makes an interest-
37 ing distance estimate. In this work, we leveraged a similar analogy between an
38 earth heap and a chromatographic elution profile. Concretely, this approach
39 is insightful since it accounts for two distinct components of what makes chro-
40 matographic elution profiles similar or not: their time separation as well as their
41 difference of shape. Let us also note that this distance has recently been applied
42 to LC-MS data, yet, to spectra rather than to chromatograms [51].

1 In general, the W1 distance between distributions \mathcal{P} and \mathcal{Q} is computed by
 2 solving Kantorovitch minimization problem, namely:

$$d_{W_1}(\mathcal{P}, \mathcal{Q}) = \inf_{\xi \in \mathcal{J}(\mathcal{P}, \mathcal{Q})} \int \|x - y\| d\xi(x, y) \quad (4)$$

3 where $\mathcal{J}(\mathcal{P}, \mathcal{Q})$ denotes all joint distributions $\xi(x, y)$ that have marginals \mathcal{P}, \mathcal{Q} .
 4 However, in the 1-dimensional discrete setting where distributions \mathcal{P} and \mathcal{Q} are
 5 replaced by chromatograms $x = (x^1, \dots, x^n)$ and $y = (y^1, \dots, y^n) \in \mathbb{R}^n$, the
 6 W1 distance boils down to a difference between empirical cumulative functions:

$$d_{W_1}(x, y) = \sum_{j=1}^n |F_x(j) - F_y(j)|, \quad (5)$$

7 where $F_x(j) = \sum_{i \leq j} \frac{x^i}{\sum_{k=1}^n x^k}$ is the j^{th} component of the cumulative distribution
 8 function of chromatogram x .

9 Kernel trick

10 Converting distances between data vectors into similarities by means of a nega-
 11 tive exponential function is a good way to derive a similarity measure endowed
 12 with the positive semi-definite (or PSD) property². This property is essential
 13 to the application of the kernel trick [59], which notably explains why kernels
 14 of the form $k(x_i, x_j) = e^{-\gamma \cdot [d_2(x_i, x_j)]^p}$, with $p = 1$ (the Laplacian kernel) or
 15 $p = 2$ (the Gaussian kernel) and with d_2 depicting the Euclidean distance are
 16 classically used.

17 Concretely, let $X = [x_1, \dots, x_N] \in \mathbb{R}^{n \times N}$ be the data matrix composed of
 18 N chromatograms. The kernel trick actually consists in using the similarity
 19 measure to implicitly map the data onto a feature space that better represents
 20 them. The mapping is deemed "implicit" as it does not require the computation
 21 of coordinates of the data point images $\Phi = [\phi(x_1), \dots, \phi(x_N)]$, where ϕ denotes
 22 the mapping function. Two conditions must be met for this trick to work: First,
 23 the algorithm must rely on similarity measures only (*i.e.* once the similarities
 24 are computed, the values of the x_i 's are not used any more). Second, the
 25 similarity measure reproduces the inner product of the feature space: $k(x, y) =$
 26 $\langle \phi(x), \phi(y) \rangle$. According to Mercer's theorem [60], any PSD similarity measure
 27 satisfies the second condition. From that point on, we refer to $K = \Phi^T \Phi =$
 28 $[k(x_i, x_j)]_{i,j=1, \dots, N}$ as the *kernel matrix*.

29 However, when using a distance like d_{W_1} , which does not derive from a
 30 norm inducing an inner product on the data space (like for instance d_2), then
 31 the PSD-ness is not guaranteed [61]. In this work, we have investigated both the
 32 Laplacian W1 and the Gaussian W1 kernels: While we exhibit a formal proof

²Positive semi-definiteness or PSD-ness, means the resulting similarity matrix will have only non-negative eigenvalues (if the eigenvalues are positive, the matrix is called positive definite or PD, see Additional File 3, Section 1).

1 of the Laplacian W1 kernel PD-ness (see Additional File 3, Section 3), we only
 2 have empirical evidence in the Gaussian case (see Additional File 3, Section 2).
 3 As in practice, both kernels lead to similar ranks in pairwise similarities, the
 4 resulting clusters only marginally differ. Owing to its popularity in life science
 5 applications, as well as to its easier tuning (interpretation and stability of the
 6 hyperparameter) the article thus focuses on the Gaussian case. Notably, as
 7 computational costs are necessarily higher with $p = 2$ than $p = 1$, the displayed
 8 runtimes are an upper bound for both cases. However, for qualitative analysis,
 9 results with $p = 1$ are also depicted in various additional files (see below).

10 Data compression

11 Nyström approximation

12 Brute force computation of a kernel matrix has a quadratic complexity, so that it
 13 does not easily scale-up. To cope for this, a classical solution is to apply Nyström
 14 approximation. This approach relies on the fast decaying property of the ker-
 15 nel spectrum (the set of kernel matrix eigenvalues): the smallest eigenvalues
 16 of the kernel matrix can safely be removed (intuitively, alike principal compo-
 17 nent analysis). Concretely, one approximates the kernel matrix $K \in \mathbb{R}^{N \times N}$ as
 18 following:

$$K \approx CW^{-1}C^\top, \quad (6)$$

19 with $C = KP \in \mathbb{R}^{N \times l}$ and $W = P^\top KP \in \mathbb{R}^{l \times l}$, where $P \in \mathbb{R}^{N \times l}$ is constructed
 20 from an $N \times N$ identity matrix where $(N - l)$ randomly selected columns are
 21 removed. The larger l , the better the approximation, but the heavier the compu-
 22 tations. Finally, according to [53], an additional rank- s truncated singular
 23 value decomposition (SVD) is of interest to increase numerical stability. This
 leads to Algorithm 1, which complexity³ is $\mathcal{O}(N \cdot n \cdot l + N \cdot l^2)$.

Algorithm 1 The rank restricted Nyström kernel approximation from [53]

- 1: **Input:** data set $X \in \mathbb{R}^{n \times N}$, similarity measure $k(\cdot, \cdot)$, Nyström sample size l , intermediate rank r , target rank s .
 - 2: Construct a random sample: $\{x_{p_1}, \dots, x_{p_l}\} \in \mathbb{R}^{n \times l}$
 - 3: Compute matrix C and W : $C = \{k(x_q, x_{p_j})\}_{\substack{q=1, \dots, N \\ j=1, \dots, l}}$, $W = \{k(x_{p_i}, x_{p_j})\}_{i,j=1, \dots, l}$.
 - 4: Perform r -truncated SVD of W : $W_r = U_r D_r U_r^\top$.
 - 5: Approximate matrix as $K \approx CW_r^{-1}C^\top = CU_r D_r^{-1} U_r^\top C^\top = RR^\top$, where $R \in \mathbb{R}^{N \times r}$.
 - 6: Perform s -truncated SVD of R : $R = U_s \Sigma_s V_s^\top$.
 - 7: **Output:** Matrix approximation $K \approx \tilde{\Phi}^\top \tilde{\Phi} = U_s \Sigma_s^2 U_s^\top$.
-

24

³As a recall, $\mathcal{O}(f(n))$ indicates that with an input data of size n , the running time will not exceed $C \cdot f(n)$ where C is a constant factor (*i.e.* independent of n).

1 It provides the following approximation of the kernel matrix: $K \approx \tilde{\Phi}^\top \tilde{\Phi}$
 2 where the matrix $\tilde{\Phi} = [\tilde{\phi}(x_1), \dots, \tilde{\phi}(x_N)]$ is obtained by applying the fea-
 3 ture mapping $\tilde{\phi}(x_i) = (\lambda_1 u_{1i}, \dots, \lambda_s u_{si})$, where λ_j and $u_{ji}, j = 1, \dots, s$ and
 4 $i = 1, \dots, N$ are the s highest eigenvalues and eigenvectors (columns of matrix
 5 U_s) of K (see Algorithm 1). Moreover, it is demonstrated in [62] that the ap-
 6 proximation accuracy is guaranteed when Nyström sample size l is on the order
 7 of \sqrt{N} . It was also shown in [53] that the target dimension s scales to $\mathcal{O}(\sqrt{l \cdot k})$,
 8 where k is the number of clusters, and the intermediate rank r is equal to $\frac{l}{2}$.

9 Random Fourier feature sketching

10 The sketching procedure of [54] is closely related to random Fourier features [63]
 11 , which seminal idea is to rely on Bochner’s theorem [64] to approximate any
 12 shift-invariant (*i.e.* $k'(x, y) = \kappa(x - y)$) PD kernel (by leveraging the fact it is
 13 a Fourier transform of some non-negative measure μ):

$$k'(x, y) = \mathbb{E}_{w \sim \mu} \left(e^{-iw^\top(x-y)} \right). \quad (7)$$

14 Elaborating on this, [54] proposed to apply a similar random Fourier map

$$\varphi(x) = \frac{1}{\sqrt{m}} \left[e^{-iw_j^\top x} \right]_{j=1}^m, \quad (8)$$

15 (where Fourier frequencies w_1, \dots, w_m are randomly sampled from some dis-
 16 tribution Ω) and to average it over all data points to approximate the data
 17 distribution itself, instead of the kernel. Concretely, applying $\varphi(\cdot)$ onto the
 18 Nyström extended data $\tilde{\Phi}$ (that is $Z = [\varphi(\tilde{\phi}(x_1)), \dots, \varphi(\tilde{\phi}(x_N))] \in \mathbb{C}^{m \times N}$), led
 19 us to computing the data sketch as:

$$SK(\tilde{\Phi}) = \frac{1}{N\sqrt{m}} \left[\sum_{i=1}^N e^{-iw_j^\top \tilde{\phi}(x_i)} \right]_{j=1}^m \in \mathbb{C}^m \quad (9)$$

20 The critical step of this data compression method lies in the frequency distri-
 21 bution estimation. It has been empirically shown in [54] that $\Omega = \mathcal{N}(0, \frac{1}{\sigma^2} \mathbf{I})$
 22 is a suitable choice for it mimicks well the fast decaying property of real life
 23 signals. Then, σ^2 can be estimated from a small data fraction using nonlinear
 24 regression. Applying this frequency distribution law allows to promote more
 25 informative sketch components and to eliminate small sketch values, which are
 26 usually related to noise. The key computational benefit of the compression is
 27 the independence between the data sketch length m and the data size N : m
 28 should be of the order of $k \cdot s$ [54], where s is the target dimension in Nyström
 29 approximation and k is the number of clusters.

1 **Cluster and centroid definitions**

2 **Cluster computations**

3 CKM (the compressive implementation of the k -means clustering presented
 4 in [56]) can be used to compute the cluster centroids from the data sketch
 5 $SK(\tilde{\Phi})$ introduced in Eq. (9). Briefly, and in contrast with classical Lloyd's
 6 algorithm, it is a greedy heuristic based on orthogonal matching pursuit, which
 7 searches for a data representation as a weighted sum of cluster centroids by
 8 minimizing the difference between corresponding sketches:

$$\|SK(\tilde{\Phi}) - \sum_{i=1}^k \alpha_i SK(c_i)\|_2^2 \quad (10)$$

9 The CKM involves two main steps summarized in Algorithm 2. First, across
 10 several iterations, it alternates between expanding the cluster centroid set with
 11 a new element, whose sketch is the most correlated to the residue; and recom-
 12 puting the centroid weights using non-negative least-squares minimization. The
 13 second step consists in the global minimization of (10) with respect to cluster
 centroids and their weights.

Algorithm 2 Compressive k-means from [56]

- 1: **Input:** data sketch $SK(\tilde{\Phi})$, frequency set w_1, \dots, w_m , the number of cen-
 centroids k , lower and upper bounds lb, ub of data $\tilde{\Phi}$.
 - 2: Initialization: $r = SK(\tilde{\Phi}), C = \emptyset$
 - 3: **for** $t \leftarrow 1$ to $2k$ **do**
 - 4: Find new centroid: $c = \arg \max_{lb \leq c \leq ub} \Re \left\langle r, \frac{SK(c)}{\|SK(c)\|} \right\rangle$
 - 5: Expand centroid set: $C = \{C, c\}$
 - 6: **if** $t > k$ **then**
 - 7: $\beta = \arg \min_{\beta \geq 0} \|SK(\tilde{\Phi}) - \sum_{i=1}^{|C|} \beta_i \frac{SK(c_i)}{\|SK(c_i)\|}\|^2$
 - 8: Choose centroids with k largest weights $C = \{c_{\beta_{i_1}}, \dots, c_{\beta_{i_k}}\}$
 - 9: **end if**
 - 10: Project to find weights: $\alpha = \arg \min_{\alpha \geq 0} \|SK(\tilde{\Phi}) - \sum_{i=1}^{|C|} \alpha_i SK(c_i)\|^2$
 - 11: Global optimization: $C, \alpha = \arg \min_{\substack{lb \leq c_i \leq ub \\ \alpha \geq 0}} \|SK(\tilde{\Phi}) - \sum_{i=1}^{|C|} \alpha_i SK(c_i)\|^2$
 - 12: Update residue: $r = SK(\tilde{\Phi}) - \sum_{i=1}^{|C|} \alpha_i SK(c_i)$
 - 13: **end for**
 - 14: **Output:** $C \in \mathbb{R}^{s \times k}$ and $\alpha_1, \dots, \alpha_k$.
-

14

1 **Cluster assignment**

2 The CKM algorithm only provides the cluster centroids and does not assign data
 3 points to clusters. Nevertheless, this can be achieved afterwards by finding the
 4 centroid which has the highest similarity value to each data point. Concretely,
 5 a cluster centroid c in the feature space can be defined using Nyström extension
 6 as follows:

$$c \approx \tilde{\phi}(y) = \Sigma_s^{-1} U_s^T k_c \quad (11)$$

7 where y is a cluster centroid in the input (chromatograms) space, and where
 8 $k_c = [k(x_1, y), \dots, k(x_N, y)]$ is an unknown vector of similarities between y and
 9 all given chromatograms. The columns of matrix U_s contain s eigenvectors of
 10 K corresponding to its s highest eigenvalues (the diagonal matrix Σ_s). The
 11 estimation of k_c can be achieved by minimizing the difference between c and
 12 $\tilde{\phi}(y)$:

$$\min_{y \in \mathbb{R}^n} \left\| \Sigma_s^{-1} U_s^T k_c - \frac{c}{\|c\|} \right\|^2 \quad (12)$$

13 The importance of the normalization term in (12) has been highlighted in [65]
 14 as an energy-preserving term to balance Nyström approximation. The solution
 15 of (12) can be found using the Moore-Penrose pseudo-inverse:

$$k_c \approx U_s \Sigma_s \frac{c}{\|c\|} \approx \tilde{\Phi}^T \frac{c}{\|c\|}. \quad (13)$$

16 Finally, the chromatographic profile x_i , $i = 1, \dots, N$ is associated to cluster j
 17 if

$$c_j = \arg \max_{c \in \{c_1, \dots, c_k\}} \left\langle \tilde{\phi}(x_i), \frac{c}{\|c\|} \right\rangle \quad (14)$$

18 The most important CKM feature is its constant execution time regardless
 19 of the data size. However, its computational complexity grows cubically with
 20 the number of clusters, so that it is not realistic to process LC-MS data where
 21 tens of thousands of clusters are classically expected. To cope for this, a divi-
 22 sive hierarchical scheme can be instrumental: Starting from a small number of
 23 clusters, one iteratively splits each cluster into k sub-clusters until a sufficiently
 24 large number of clusters k_{total} is achieved. However, this strategy requires, for
 25 each independent call of the clustering algorithm, an update of the data sketch
 26 as well as a complete assignment to clusters. Thus, to practically improve its
 27 computational efficiency, we leveraged the expected decrease of the cluster size
 28 at each iteration to optimize the code, and we decided to compute all the data
 29 sketches from the same frequency samples, either on the entire dataset (at first
 30 step) or on the cluster to be re-clustered (at the following iterations). Finally, it
 31 appeared these repetitive computations of the cluster sketches and assignments
 32 did not hamper the efficiency of the whole process.

33 **Pre-image computation**

34 The combination of Nyström approximation and of random Fourier features
 35 leads to an additional difficulty: To recover the signal of each consensus elution

1 profile, it is necessary to compute its reverse mapping from the feature space
 2 back to the input space. This is referred to as a *pre-image* problem and it is ill-
 3 posed: only an approximation of the cluster centroids in the input space can be
 4 obtained. The conventional fixed point iteration method [66] cannot be applied
 5 due to the use of the W1 distance. Similarly, the reconstruction of a consensus
 6 chromatogram as the mean of the cluster chromatograms is not adapted, due to
 7 large scale non-linearities between the input and feature spaces, as illustrated
 8 in Figure 1C.

9 To correct for this, we decided to compute a local (*i.e.* small-scale) mean
 10 by considering only a subset of the closest chromatograms. To determine the
 11 cluster centroid neighbourhood $\mathcal{N}(c)$, we proceeded similarly to the cluster as-
 12 signment step, by choosing the chromatograms in the cluster $\mathcal{J}(c)$ with the
 13 highest similarities to the cluster centroid:

$$\mathcal{N}(c) = \{x_1, \dots, x_q\} \subset \mathcal{J}(c) \mid k(c, x) > k(c, y) \quad \forall x \in \mathcal{N}(c), y \in \mathcal{J}(c) \setminus \mathcal{N}(c), \quad (15)$$

14 where similarities $k(c, \cdot)$ were estimated using Eq. (13). Concretely, $\mathcal{N}(c)$ was
 15 defined by selecting the q closest neighbors (so that $q = |\mathcal{N}(c)|$). The tuning
 16 of parameter q is discussed with that of other parameters in the Parameter
 17 tuning section.

18 Performance metrics

19 For experiments annotated with a ground truth (like UPS2GT dataset), clus-
 20 tering accuracy can be evaluated with the Rand index (RI). The Rand index
 21 measures the percentage of correctly clustered pairs of signals over the total
 22 number of pairs. Let us denote as $U = \{U_1, \dots, U_k\}$ the obtained clusters and
 23 as $V = \{V_1, \dots, V_q\}$ the ground truth clusters. A pair of signals is considered
 24 as correctly clustered: *true positive (TP)* or *true negative (TN)*, if signals are
 25 assigned to the same cluster in U and V or on the contrary, to different clusters
 26 in U and V . A pair of signals is called *false positive (FP)* (resp. *false negative*
 27 *(FN)*), if signals are grouped in U (resp. V) but not in V (resp. U). Then, the
 28 Rand index is given by:

$$\text{RI} = \frac{TP + TN}{TP + TN + FP + FN} \quad (16)$$

29 The maximum value of the Rand index is 1 (perfect match with the ground
 30 truth). Additionally, it is possible to evaluate how often different chromatograms
 31 are grouped in the same cluster; and how often similar chromatograms were as-
 32 signed to different clusters. To do so, one classically relies on the Precision and
 33 Recall metrics, respectively:

$$\text{Precision} = \frac{TP}{TP + FP} \quad \text{Recall} = \frac{TP}{TP + FN} \quad (17)$$

34 For datasets without ground truth annotation (like both Ecoli datasets), it
 35 is possible to rely on the Davies - Bouldin (DB) index. Let us denote as $\mathcal{J}(c_j)$

1 the j^{th} cluster with the cluster centroid c_j , and as $\{\mathcal{J}(c_1), \dots, \mathcal{J}(c_k)\}$ the set
 2 of obtained clusters. The within cluster distance reads:

$$S_j = \frac{1}{|\mathcal{J}(c_j)|} \sum_{x_i \in \mathcal{J}(c_j)} d_{W_1}(x_i, c_j) \quad (18)$$

3 The DB index is defined through the ratio of the within cluster distances to the
 4 between cluster distance $d_{W_1}(c_i, c_j)$:

$$DB = \frac{1}{k} \sum_{i=1}^k \max_{i \neq j} \frac{S_i + S_j}{d_{W_1}(c_i, c_j)}, \quad (19)$$

5 It should be noted that the distance metric in the DB index and in the clus-
 6 tering algorithm must be the same, in our case the W1 distance in the original
 7 space. Moreover, the smaller the DB index, the better the clustering (as a good
 8 clustering minimizes cluster overlaps).

9 Finally, the computational load can easily be approximated by the recorded
 10 execution time, *i.e.* the difference between the end and start times, both of
 11 which being accessible in R with the `Sys.time()` function. For sake of brevity,
 12 execution times are reported for the Gaussian W2 kernel only, as Laplacian
 13 similarities are necessarily faster to compute (no squared distance to evaluate).

14 Results

15 Objectives of the experimental assessment

16 Many independent elements deserve evaluations: The first one is the practi-
 17 cal interest of W1 distance in the context of LC-MS data. The second one is
 18 the computational load of our complete algorithm in function of the parame-
 19 ter tuning (on the one hand, an efficient compression technique is used; on the
 20 other hand, one targets the clustering of raw data into a high number of clus-
 21 ters, making its efficiency a challenge). The third one is the clustering result
 22 itself. However, a classical evaluation of the clustering performances will be of
 23 little interest: In fact, all k -means related algorithms (including their kernel-
 24 ized versions) have been extensively studied [44], so that their strengths and
 25 weaknesses are now well-documented. For instance, k -means optimizers can
 26 easily be trapped into local minima and cannot naturally deal with outliers,
 27 which are both significant drawbacks; however, they scale up well to very-high
 28 dimensional data, which definitely is an asset for LC-MS applications. In con-
 29 trast, highlighting the differences of our approach with respect to linkage-based
 30 agglomerative clustering and showing that despite noticeable differences, one
 31 obtains clusters which are meaningful, is of real practical interest to computa-
 32 tional mass spectrometry experts.

33 As reported in the Background section, comparisons with Xnet is mandatory.
 34 However, considering the reported specificities (trace extraction preprocessing,

1 envelope assumption simplification, *etc.*), comparing Xnet and CHICKN work-
2 flows may appear as somewhat arbitrary. To cope for this, we have made the
3 following choices: First, we have focused on the core of each algorithm, as repre-
4 sented in Figure 2A. Second, we have adapted the UPS2GT and Ecoli datasets
5 to be processed by each algorithm: The UPS2GT data are already formatted
6 into a CSV file meeting Xnet requirements. To construct a data matrix suit-
7 able to CHICKN from the UPS2GT data, we simply loaded the data points
8 according to their retention time and trace labels in the matrix columns (simi-
9 larly to Xnet, we excluded point with trace indices -1 and 0, as assumed to be
10 noise). This led to a data matrix containing 57,140 columns and 6,616 rows.
11 Conversely, to build the CSV files from Ecoli datasets, we stored any non-zero
12 entry of the data matrix in a row, the column index being used in place of the
13 trace labels.

14 Wasserstein distance validation

15 W1 distance was proposed to discriminate between signals that represent dif-
16 ferent elution profiles. To assess this choice, we compared it with two distances
17 amongst the most widely used in mass spectrometry signal processing: The first
18 one is the classical Euclidean distance. The second one is the peak retention
19 time difference (or ΔRT): It corresponds to the difference between the time
20 stamps at which each signal reaches its highest intensity value. Based on the
21 Ecoli-FMS dataset (which provides the finest temporal sampling), we examined
22 two situations presented in Figure 3: In the first one, we selected 3 signals
23 with different shapes, that we precisely aligned so that their pairwise ΔRT
24 was zero; in other words, only the shape difference makes it possible to discriminate
25 them. Conversely, in the second situation, an elution profile was translated to
26 mimic a case where only the ΔRT was meaningful. In both situations, the sec-
27 ond chromatogram (chr2) stands as an in-between the first (chr1) and the third
28 chromatogram (chr3). As illustrated by the distance ratios given in the tables
29 embedded in Figure 3, both the Euclidean and the ΔRT distances are meaning-
30 ful in one case: The Euclidean distance captures the shape information, while
31 ΔRT captures the time translation effect. However, none of these classically
32 used distances is able to capture both the shape and the translation simultane-
33 ously. On the contrary, W1 distance is efficient on both situations, making it a
34 suitable distance to construct a similarity measure adapted to LC-MS data.

35 Parameter tuning

36 Unlike Xnet, CHICKN is governed by eight parameters. Four of them are
37 involved in the data compression: Nyström sample size (l), target rank (s),
38 kernel parameter (γ) and sketch size (m). Three parameters are involved in the
39 hierarchical clustering: number of clusters at each iteration of the hierarchical
40 clustering (k), upper bound of the total number of expected clusters (k_{total}) and
41 maximum number of levels in the hierarchy (T). The remaining parameter is the
42 neighbourhood size in the consensus chromatogram computation (q). However,

1 all parameters except γ and q are interrelated (see the Data compression section
 2 as well as [62, 53]) and can be expressed through k , k_{total} and N (the dataset
 3 size) as follows:

$$\begin{aligned}
 l &\approx \sqrt{N}, \\
 s &\approx \sqrt{k} \cdot N^{1/4}, \\
 m &\approx k^{3/2} \cdot N^{1/4}, \\
 T &= \lfloor \log(k_{\text{total}}, k) \rfloor.
 \end{aligned}
 \tag{20}$$

4 These theoretical results can nonetheless be discussed. Notably, tuning the
 5 sketch size m to a larger value may be of interest if contrarily to our case, the
 6 computational efficiency is not the only targeted goal. Thus, we have performed
 7 complementary investigation to relate the clustering performance (in terms of
 8 DB index) to the sketch size (see Additional File 4, leftmost figure). Oddly
 9 enough, it appears the DB index increases (*i.e.* the performances deteriorates)
 10 when the sketch size increases (leading to a more refined representation of the
 11 data). However, it appears to be an indirect consequence: when increasing m ,
 12 more differences between the signals are represented, making it possible to define
 13 a larger number of smaller clusters (see Additional File 4, rightmost figure).

14 Finally, four parameters remain (γ , q , k and k_{total}). Concretely, we tuned
 15 the kernel parameter γ as an average of the power of p distances to the ν nearest
 16 neighbors for all chromatographic profiles:

$$\gamma = \frac{1}{N \cdot \nu} \sum_{i=1}^N \sum_{j=1}^{\nu} [d_{W_1}(x_i, x_{i_j})]^p,
 \tag{21}$$

17 where $x_{i_1}, \dots, x_{i_\nu}$ are ν neighbors of x_i (selected among the l points of the
 18 Nyström sample) and $p \in \{1, 2\}$ depending on the kernel type. Practically,
 19 we observed that tuning ν to 32 guaranteed each data point to be sufficiently
 20 connected to the rest of the dataset, as advised in [37]. Moreover, we observed
 21 that γ was rather stable with respect to ν , for both Laplacian W1 and Gaussian
 22 W1 kernels. However, as expected, the stability is higher with the latter than
 23 with the former (see Additional File 5).

24 For q (in the consensus chromatogram computation) we observed that the
 25 shape cluster problem (see Pre-image computation section) could only occur
 26 with significantly large clusters (few tenth of elements). Thus, as preliminary
 27 stability analysis indicated us that the consensus chromatogram shapes were
 28 preserved across various values of ν (see Additional File 6), we decided to bound
 29 q with ν and to set $q = \min(\nu, \text{cluster size})$.

30 A known drawback of k -means objective function is the requirement to set
 31 the maximum number of expected clusters (knowing some clusters can remain
 32 empty). In our case, this is achieved by tuning k and k_{total} . Yet, it should
 33 be noted that increasing k leads to decreasing T for a fixed value of k_{total} so
 34 that a trade-off between T and k must be sought. With this respect, we have
 35 evaluated different scenarios with $k = 2, 4, 8$ and 16. CHICKN execution times
 36 (excluding the data compression step, which remains constant whatever the

1 various scenario) on the smallest (UPS2GT) and largest (Ecoli-FMS) datasets
 2 are depicted in Additional File 7. This experiment pointed out the importance
 3 of tuning k to a small enough value, which is coherent with the observation that
 4 the original CKM algorithm does not scale up well with the number of clusters.
 5 Practically, working with $k = 2$ or 4 appeared to be the most efficient.

6 In the case of UPS2GT, the expected number of isotopic envelopes is known
 7 (*i.e.* 14,076). Thus, it is easy to tune k_{total} accordingly (*i.e.* $2^{14} = 4^7 = 16,384$).
 8 However, knowing that CHICKN does not rely on the envelope assumption
 9 simplification, it can be expected to find a much lower number of clusters:
 10 broadly, all the isotopic envelopes corresponding to different charge states of
 11 a same peptide can be expected to cluster together. Therefore, it also makes
 12 sense to tune k_{total} to $4^5 = 1,024$; *i.e.* close enough from the expected number
 13 of identifiable peptides in the sample (around 700, according to [23]).

14 Tuning k_{total} for any real life data (*i.e.* unlabeled) is much more complicated.
 15 However, the *Escherichia Coli* sample is well studied, and based on prior biolog-
 16 ical/analytical knowledge, 15,000 different peptides can be expected, broadly.
 17 Consequently, for both Ecoli datasets, $k_{\text{total}} = 16,384$ seems reasonable. Fi-
 18 nally, even though it is not as sensible from a biological viewpoint, we have
 19 decided to also consider $k_{\text{total}} = 4^6 = 4,096$, which provides an even ground for
 20 computational load comparisons (see next section for details).

21 To summarize, three different ways to tune k_{total} are insightful: 1,024 for the
 22 UPS2GT dataset only (as it matches the number of expected peptides); 4,096
 23 on all datasets (for computational benchmarks); and 16,384 on all datasets
 24 (number of isotopic envelopes in UPS2GT and number of expected peptides in
 25 Ecoli datasets).

26 Finally, we fixed the remaining parameter values using the formulas in Eq.
 27 (20), as summarized in Table 1.

Table 1: Summary of the different combinations of parameter tuning.

Dataset	γ	l	s	m	k	T		
						$k_{\text{total}} = 1,024$	$k_{\text{total}} = 4,096$	$k_{\text{total}} = 16,384$
UPS2GT	5.96e-06	240	22	44	2	10	12	14
	6.9e-06	240	31	124	4	5	6	7
Ecoli-DIA	9.06e-06	432	30	60	2	-	12	14
	9.27e-06	432	42	168	4	-	6	7
Ecoli-FMS	7.07e-07	863	42	84	2	-	12	14
	7.03e-07	863	59	236	4	-	6	7

28 Computational load

29 We have compared the execution times of CHICKN and Xnet cores (see Fig-
 30 ure 2A). Previously reported comparisons showed us that CHICKN execution
 31 time largely depends on k . However, it only has a sub-linear complexity with
 32 respect to k_{total} . : As illustrated in Additional File 8, multiplying k_{total} by 4

1 only results in a threefold (resp. twofold) increase in the CHICKN run-time for
2 the Ecoli-FMS (resp. UPS2GT) dataset. As reducing k_{total} to limit the execu-
3 tion time will therefore be of little interest, experiments hereafter reported only
4 focused on the influence of k . Despite CHICKN being more efficient when run
5 with $k = 2$ and 4 (see Parameter tuning section), we also included comparisons
6 with $k = 8$ and 16 to investigate the consequences of sub-optimal parame-
7 ter tuning. The corresponding tests are referred to as CHICKN2, CHICKN4,
8 CHICKN8 and CHICKN16. Therefore, to rely on an even basis for comparisons,
9 we focused on $k_{\text{total}} = 4,096$: it is a power of 16, contrarily to 1,024 and 16,384
10 (which are even not a power of 8).

11 Since CHICKN algorithm embeds a compressive k -means algorithm which
12 may converge towards different local minima depending on the stochasticity of
13 several steps, each scenario was repeated 10 times and the average execution
14 time was reported. In contrast, Xnet being deterministic, it was executed once.
15 In [17], Xnet exhibits impressive computational times on pre-processed and
16 adequately formatted data. However, raw LC-MS data stored in a matrix format
17 are more cumbersome. Thus, our first experiment was to compare the efficiency
18 of Xnet and of CHICKN on the Ecoli-DIA dataset, using a laptop machine
19 with the following characteristics: HP Pavilion g6 Notebook PC with Intel(R)
20 Core(TM) i5-3230M CPU @ 2.60GHz, 8 Gb of RAM, 4 cores, running under
21 Ubuntu 18.04.4 LTS OS. Xnet produced an "out-of-memory" error when trying
22 to cluster more than 10,000 columns (*i.e.* 5% of the Ecoli-DIA dataset) in
23 a single batch. This is why Figure 2B compares the computational time of
24 CHICKN2, CHICKN4, CHICKN8 and of CHICKN16 on the entire Ecoli-DIA
25 dataset to that of Xnet on only 5% of the same dataset. On this figure, different
26 colors are used to discriminate between the clustering step *per se* and CHICKN
27 preliminary data compression step. Let us note that the compression step is time
28 consuming, however, it also includes the computations of all the W1 similarities.
29 This as-a-matter-of-factly illustrates the computational cost of relying on more
30 elaborated metrics to capture the semantics of data as complex as LC-MS ones.
31 Except for CHICKN16, which has already been pointed as suboptimal, CHICKN
32 is always faster for a dataset 20 times larger.

33 This first experiment clearly showed CHICKN could be used on a simple
34 laptop, even with large datasets, in long but acceptable times (half an hour to
35 two hours, broadly). Then, to reduce the execution times of our multiple experi-
36 ments, but also to allow Xnet working on a larger dataset, we moved to a larger
37 station using 10 cores of an Intel Xeon CPU E5-2470 v2 @ 2.40GHz, 94 GB
38 of RAM and running with CentOS Linux release 7.4.1708. As depicted in Fig-
39 ure 2C, on such a machine, CHICKN was able to process Ecoli-FMS within 5h30
40 (most of them being necessary to perform the preliminary compression), despite
41 its huge size. On the contrary, with the same machine, Xnet only processed 10%
42 of it in a comparable time (almost 8 hours). Moreover, larger fractions of the
43 dataset were not processable, as leading to memory failure.

44 To explain this discrepancy, we noticed that Xnet spent a considerable time
45 to construct the preliminary network. The nature of Ecoli data (raw data
46 without any trace pre-processing and recorded with the highly resolved *profile*

1 mode, see Materials section) contrasts with that of UPS2GT, on which Xnet is
2 really efficient. As it appears on Figure 2D, CHICKN is clearly not as fast as
3 Xnet to process UPS2GT: The Xnet analysis took less than 40 seconds, while
4 CHICKN computation times varied from 2 to 7 minutes depending on values of
5 parameter k (from 2 to 16).

6 As a whole, these experiments illustrate the utmost importance of prior
7 preprocessing methods when studying LC-MS data. In this context, algorithms
8 working on raw data, such as CHICKN, are real assets.

9 Cluster evaluation

10 Figure 4 reports the Rand index, Precision and Recall (UPS2GT dataset) as well
11 as the DB index (Ecoli datasets) with different clustering strategies: CHICKN2
12 and CHICKN4 (with $k_{\text{total}} \in \{1,024 ; 4,096 ; 16,384\}$ and with $p = 2$), as well
13 as Xnet (on UPS2GT only, for computational reasons). A similar figure for
14 $p = 1$ is available in Additional File 9.

15 First, it can be noted that the Rand index is hardly informative (Figure 4A):
16 All clustering methods exhibit an index of almost 1, and it is necessary to go
17 three (and sometimes four) decimals to notice a difference. Such high values
18 are a direct consequence of the huge number of expected clusters in UPS2GT
19 datasets, which comes with an excessively large number of true negative pairs
20 (almost 99 % of all possible pairs). In this context, the Rand index obtained
21 with "only" 1,024 expected clusters is particularly highlighting: Despite 16 times
22 less clusters, it achieves an equivalent index. This indicates that, relatively, the
23 provided clustering is probably of better quality.

24 However, contrarily to the Rand index, Precision and Recall are informative
25 to compare with Xnet, as the true negative pair count does not level the scores.
26 With this regard, it clearly appears on Figures 4B that the Precision is in-
27 comparably better with Xnet. Although foreseeable (ground truth with 14,076
28 envelopes whereas CHICKN sought a thousand of peptides), this requires a
29 deeper analysis: Concretely, Xnet tends to over-cluster (which artificially im-
30 proves the Precision index), as it provided 17,153 clusters covering 93% of the
31 dataset (7% of the elution profiles are excluded by Xnet) where the ground truth
32 labels proposed only 14,076 of them (on 100% of the dataset). In addition, Xnet
33 priors were trained on the same UPS2GT dataset as for evaluation, so that high
34 performance are expectable. With this regard, it is particularly noteworthy that
35 the Recall (Figures 4C) varies the other way around. Concretely, it is best for
36 CHICKN4 with $k_{\text{total}} = 1,024$ despite this number being completely different
37 from the one derived from the ground truth. In addition to be in line with our
38 observations on the Rand index, this concurs with the peptide-level knowledge
39 of the dataset: CHICKN was supposed to group together differently charged
40 peptides, which it did (see Additional Files 10 and 14 as well as Discussions
41 below), as it provided only 510 (CHICKN4)/ 740 (CHICKN2) clusters on the
42 entire UPS2GT dataset, hereby leaving 300 to 500 empty clusters⁴; and leading

⁴More generally, the capability of CHICKN to adapt the cluster sizes to the data distribu-

1 to a number of clusters in line with the expected number of peptides in the
2 sample. Overall, the differences between Xnet and CHICKN on UPS2GT seem
3 to be more related to the difference of objectives (finding isotopics envelopes
4 *vs.* finding peptide-related clusters), as already discussed. Interestingly, this
5 interpretation is confirmed by the Ecoli dataset experiments.

6 In absence of ground truth for both Ecoli datasets, we chose the tuning mini-
7 mizing the DB index (see Figure 4D and 4E): $k_{\text{total}} = 16,384$ for Ecoli-FMS and
8 for Ecoli-DIA. With such a tuning, we obtained around 11,600 (resp. around
9 9,400) non-empty clusters for Ecoli-FMS (resp. Ecoli-DIA). This number is
10 obviously lower than the expected number of identifiable peptides (between 15
11 and 20 thousands), however under-clustering was clearly supported by empirical
12 observations (see above, as well as Additional File 4, rightmost figure). This
13 clearly means that CHICKN could not separate too many peptides with too
14 similar elution profiles. However, this can be easily explained by the difference
15 of complexity between the UPS2GT and the Ecoli samples: while the former
16 is fairly simple (a handful of spiked proteins), the latter ones are complex real
17 life samples for which the discriminative power of the liquid chromatography is
18 clearly challenged (as illustrated in the next section). This is notably why frag-
19 mentation spectra are classically used to identify as many as 15 to 20 thousand
20 peptides. However, achieving to discriminate half of this number of peptides
21 with MS1 processing only is noticeable.

22 Finally, let us note, that, in general, relying on $k = 4$ provided slightly better
23 scores. We assume that $k = 4$ was a trade-off between cluster diversity ($k > 4$)
24 and computational efficiency ($k = 2$), as discussed above.

25 Discussions

26 Cluster interpretability

27 Beyond evaluation metrics, it is insightful to compare algorithms according to
28 the interpretability of the clusters they can provide. Figure 5 represents differ-
29 ent elution profiles from UPS2GT (their shape as well as their m/z position) in
30 the context of the clusters they fall into, according to CHICKN and Xnet. The
31 envelope assumption simplification clearly appears: As expected, Xnet splits
32 into different clusters elution profiles that are arguably similar for the reason
33 they have too different m/z values. In contrast, CHICKN promotes the inner
34 coherency of clusters as it aggregates related Xnet clusters together. Notably,
35 Additional Files 10 and 14 show a subset of 12 clusters provided by CHICKN,
36 each gathering at least 2 differently charged ions from a same peptide (all of
37 them being identified and manually validated with the associated MS2 spec-
38 tra). Interestingly, the multiple isotopes of each ion also appear to be grouped,
39 as illustrated by the manifold of profile co-clustered with each ion. Moreover,
40 a refine analysis of CHICKN clusters shows that, globally, they contain similar
41 chromatograms, which is coherent both with the clustering metrics provided

tion is illustrated on Additional File 11.

1 above, and with the expected behavior of the W1 kernel. However, some clusters
2 also contain noise signals, as for examples, the first two lines of Figure 5.
3 Although undesirable, this is a direct consequence of (i) the grouping capabilities
4 of CHICKN, which captures similarities between slightly different but
5 largely overlapping signals (third line); and (ii) the possibility to run CHICKN
6 on raw data, which also contains many spurious signals that need be spread
7 across various meaningful clusters.

8 Similar conclusions regarding CHICKN behavior can be derived from the
9 Ecoli datasets (let us focus on the Ecoli-FMS one, as it displays elution profile
10 signals with higher sampling resolution, due to the Full-MS acquisition). The
11 majority of clusters (Figure 6 for the Gaussian W1 kernel and Additional File 12,
12 for the Laplacian W1 one) containing high intensity signals depicts meaningful
13 consensus chromatograms, as well as similar profiles even though corresponding
14 to different m/z values. However, we observed that some clusters could be separated
15 into several sub-clusters to improve readability (see Additional File 13).
16 It could intuitively be interpreted as the necessity to increase k_{total} . However,
17 two observations goes against this: First, from a signal viewpoint, as the phenomenon
18 mainly impacts lower intensity profiles, it also highlights the difficulty
19 of finding consensus patterns near the noise level, which equally affects most
20 of the clustering algorithms. In this context, over-clustering is usually not considered
21 a viable solution. Second, from an analytical viewpoint, the clustering
22 algorithm cannot be expected to separate beyond the chromatographic capabilities
23 (as in Additional File 13, where few different profiles have too important
24 overlap to expect discrimination).

25 Finally, it is worthy focusing on consensus chromatograms: interestingly
26 enough, most of those observed in Figure 6 and in Additional File 12 have
27 meaningful shapes that are not deteriorated by the presence of noisy signals in
28 the cluster, which can be interpreted as a positive consequence of our method
29 to compute the cluster centroids pre-image based on a restricted neighborhood
30 (see Pre-image computation section).

31 **Implementation and code availability**

32 CHICKN algorithm was implemented in R. The W1 distance computations and
33 the gradient descent were accelerated using C and interfaced with R thanks
34 to Rcpp. The data compression procedure and the hierarchical strategy were
35 parallelized with RcppParallel, foreach and doParallel. To access and manipulate
36 large data matrices, we relied on the File-backed Big Matrix class of the
37 bigstatsr package [67]. A File-backed matrix allows to overcome the memory
38 limitation by storing the data on the disk, using a binary memory-mapped
39 file. However, bigstatsr is only available under Linux OS, leading to a similar
40 restriction for CHICKN.

41 For practitioners, the proposed algorithm is available through an R package,
42 available on Gitlab [68], as well as on the CRAN [69].

1 **Conclusion**

2 We have presented two complementary contributions to the cluster analysis of
3 LC-MS data. First, we have proposed a unique combination of hierarchical
4 strategy, of Nyström approximation and of random Fourier features based com-
5 pression technique to scale up the kernel k -means clustering to the large size,
6 the large dimensionality and the large number of expected clusters of LC-MS
7 data. Second, we have proposed to rely on the optimal transport framework
8 (Wasserstein-1 distance) to define a similarity measure and we have shown it is
9 insightful to capture the semantics of elution profiles in LC-MS data. On a more
10 theoretical front, we have established the Wasserstein-1 distance could lead to a
11 positive-definite Laplacian kernel, and exhibit a path for further investigations
12 about a Gaussian one.

13 We have demonstrated these contributions could help extracting other struc-
14 tures than isotopic envelopes, even on multiplexed data acquired with Data
15 Independent Acquisition protocol. However, the experimental assessment of
16 these contributions is difficult to interpret. On the one hand, when compared
17 to the canonical application of isotopic envelope extraction, CHICKN does not
18 outperform the state-of-the-art algorithm (better Recall and worse Precision,
19 as it tends to under-cluster rather than over-cluster). However, it provides an
20 important advantage: it can be run on raw data and does not require costly pre-
21 processing. As for an application-independent evaluation, it clearly appears that
22 CHICKN is able to extract patterns from the data which are not accessible to
23 linkage-based algorithms. Put together, we interpret this as following: Although
24 cluster analysis has made important progresses in the theoretical front over the
25 past 50 years, processing LC-MS data remains a challenge which requires re-
26 search efforts. It is still necessary to propose complementary and differently
27 principled algorithms that will help make LC-MS practitioners extract the best
28 from their data. In this context, new kernels could be defined; and numerous
29 state-of-the-art clustering algorithms recently developed in the machine learning
30 community could advantageously be applied to LC-MS data.

31 **Declarations**

32 **Ethics approval and consent to participate**

33 Not applicable

34 **Consent for publication**

35 Not applicable

36 **Availability of data and materials**

37 The UPS2GT dataset supporting the conclusions of this article is available in
38 the Github, <https://github.com/optimusmoose/ups2GT> [47].

1 The Ecoli DIA and Ecoli FMS raw data (generated and analyzed for the
2 current study) are not publicly available due to their too large size. However
3 they are available from the corresponding author upon reasonable request. To
4 reproduce all experiments described in the article, the preprocessed Ecoli datasets
5 in the file-backed matrix format can also be provided upon request.

6 **Funding**

7 This work was supported by grants from the French National Research Agency:
8 ProFI project (ANR-10-INBS-08), GRAL project (ANR-10-LABX-49-01), DATA@UGA
9 and SYMER projects (ANR-15-IDEX-02) and MIAI @ Grenoble Alpes (ANR-
10 19-P3IA-0003). Grants ANR-10-INBS-08 and ANR-10-LABX-49-01 contributed
11 to wet-lab equipment, including mass spectrometer; Grant ANR-15-IDEX-02
12 contributed to human resources; Grant ANR-19-P3IA-0003 supported other ex-
13 penses.

14 **Competing interests**

15 The authors declare that they have no competing interests.

16 **Author's contributions**

17 OP designed the method, implemented the R package, carried out the compu-
18 tational experiments, analysed the results and drafted the manuscript. RG im-
19 plemented Nyström approximation method and the Ecoli dataset construction
20 routines. TF designed the preprocessing steps, performed preliminary imple-
21 mentations and advised on some mathematical issues. AK and AMH produced
22 the Ecoli datasets (LC-MS analysis design and production, preliminary bioin-
23 formatics processing). AMH wrote the wet-lab analysis section. TB designed
24 the method, directed the work, participated to the result analysis and drafted
25 the manuscript. All authors proofread the manuscript and approved its final
26 version.

27 **Acknowledgements**

28 The authors thank Virginie Brun, Yohann Couté and Christophe Bruley for
29 supports and fruitful discussions.

30 **List of abbreviations**

- 31 • **CHICKN**, Chromatogram Hierarchy Compressive K-means with Nys-
32 trom approximation;
- 33 • **CKM**, Compressive k-means;
- 34 • **DB**, Davies-Bouldin index;

- 1 • **Δ RT**, peak Retention Time difference;
- 2 • **DIA**, Data Independent Acquisition;
- 3 • **Ecoli-DIA**, Dataset resulting from the analysis of an Escherichia coli
- 4 sample in DIA mode;
- 5 • **Ecoli-FMS**, Dataset resulting from the analysis of an Escherichia coli
- 6 sample in FMS mode;
- 7 • **FMS**, Full mass spectrum (only MS1s are recorded);
- 8 • **FN**, False Negative;
- 9 • **FP**, False Positive;
- 10 • **LC-MS**, liquid chromatography and mass spectrometry;
- 11 • **MS1**, peptide (or precursor) mass spectrum;
- 12 • **MS2**, peptide fragment mass spectrum;
- 13 • **m/z**, mass to charge ration;
- 14 • **PD**, positive definite;
- 15 • **PSD**, positive semi-definite;
- 16 • **RI**, Rand index;
- 17 • **RT**, Retention Rime;
- 18 • **SVD**, Singular Value Decomposition;
- 19 • **TN**, True Negative;
- 20 • **TP**, True Positive;
- 21 • **UPS2GT**, dataset resulting from the analysis of the Proteomics Dynamic
- 22 Range Standard (UPS2) and manually annotated (ground truth);
- 23 • **W1**, Wasserstein-1 distance;

24 **References**

- 25 [1] Teleman J, Dowsey AW, Gonzalez-Galarza FF, Perkins S, Pratt B, Röst
- 26 HL, et al. Numerical compression schemes for proteomics mass spectrom-
- 27 etry data. *Molecular and Cellular Proteomics*. 2014;13(6):1537–1542.
- 28 [2] Klaus B, Strimmer K. Signal identification for rare and weak features:
- 29 Higher criticism or false discovery rates? *Biostatistics*. 2013;14(1):129–
- 30 143.

- 1 [3] Tabb DL, MacCoss MJ, Wu CC, Anderson SD, Yates JR. Similarity among
2 tandem mass spectra from proteomic experiments: Detection, significance,
3 and utility. *Analytical Chemistry*. 2003;75(10):2470–2477.
- 4 [4] Tabb DL, Thompson MR, Khalsa-Moyers G, VerBerkmoes NC, McDonald
5 WH. MS2Grouper: Group assessment and synthetic replacement of dupli-
6 cate proteomic tandem mass spectra. *Journal of the American Society for
7 Mass Spectrometry*. 2005;16(8):1250–1261.
- 8 [5] Beer I, Barnea E, Ziv T, Admon A. Improving large-scale proteomics by
9 clustering of mass spectrometry data. *Proteomics*. 2004;4(4):950–960.
- 10 [6] Flikka K, Meukens J, Helsens K, Vandekerckhove J, Eidhammer I, Gevaert
11 K, et al. Implementation and application of a versatile clustering tool for
12 tandem mass spectrometry data. *Proteomics*. 2007;7(18):3245–3258.
- 13 [7] Frank AM, Bandeira N, Shen Z, Tanner S, Briggs SP, Smith RD, et al.
14 Clustering millions of tandem mass spectra. *Journal of Proteome Research*.
15 2008;7(1):113–122.
- 16 [8] Frank AM, Monroe ME, Shah AR, Carver JJ, Bandeira N, Moore RJ, et al.
17 Spectral archives: Extending spectral libraries to analyze both identified
18 and unidentified spectra. *Nature Methods*. 2011;8(7):587–594.
- 19 [9] Griss J, Foster JM, Hermjakob H, Vizcaíno JA. PRIDE Cluster: Building
20 a consensus of proteomics data. *Nature Methods*. 2013;10(2):95–96.
- 21 [10] Griss J, Perez-Riverol Y, Lewis S, Tabb DL, Dienes JA, Del-Toro N, et al.
22 Recognizing millions of consistently unidentified spectra across hundreds of
23 shotgun proteomics datasets. *Nature Methods*. 2016;13(8):651–656.
- 24 [11] Falkner JA, Falkner JW, Yocum AK, Andrews PC. A spectral cluster-
25 ing approach to MS/MS identification of post-translational modifications.
26 *Journal of Proteome Research*. 2008;7(11):4614–4622.
- 27 [12] Saeed F, Hoffert JD, Knepper MA. CAMS-RS: Clustering algorithm for
28 large-scale mass spectrometry data using restricted search space and in-
29 telligent random sampling. *IEEE/ACM Transactions on Computational
30 Biology and Bioinformatics*. 2014;11(1):128–141.
- 31 [13] The M, Käll L. MaRaCluster: A Fragment Rarity Metric for Clustering
32 Fragment Spectra in Shotgun Proteomics. *Journal of Proteome Research*.
33 2016;15(3):713–720.
- 34 [14] Griss J, Perez-Riverol Y, The M, Käll L, Vizcaíno JA. Response to "com-
35 parison and Evaluation of Clustering Algorithms for Tandem Mass Spec-
36 tra". *Journal of Proteome Research*. 2018;17(5):1993–1996.
- 37 [15] Wang L, Li S, Tang H. MsCRUSH: Fast Tandem Mass Spectral Clus-
38 tering Using Locality Sensitive Hashing. *Journal of Proteome Research*.
39 2019;18(1):147–158.

- 1 [16] Perez-Riverol Y, Vizcaíno JA, Griss J. Future Prospects of Spectral Clus-
2 tering Approaches in Proteomics. *Proteomics*. 2018;18(14):1700454.
- 3 [17] Gutierrez M, Handy K, Smith R. XNet: A Bayesian Approach to Ex-
4 tracted Ion Chromatogram Clustering for Precursor Mass Spectrometry
5 Data. *Journal of Proteome Research*. 2019;18(7):2771–2778.
- 6 [18] Fischer B, Grossmann J, Roth V, Gruissem W, Baginsky S, Buhmann JM.
7 Semi-supervised LC/MS alignment for differential proteomics. *Bioinfor-*
8 *matics*. 2006;22(14):e132—e140.
- 9 [19] Houel S, Abernathy R, Renganathan K, Meyer-Arendt K, Ahn NG, Old
10 WM. Quantifying the impact of chimera MS/MS spectra on peptide iden-
11 tification in large-scale proteomics studies. *Journal of Proteome Research*.
12 2010;9(8):4152–4160.
- 13 [20] Chapman JD, Goodlett DR, Masselon CD. Multiplexed and data-
14 independent tandem mass spectrometry for global proteome profiling. *Mass*
15 *Spectrometry Reviews*. 2014;33(6):452–470.
- 16 [21] Peckner R, Myers SA, Jacome ASV, Egertson JD, Abelin JG, MacCoss
17 MJ, et al. Specter: Linear deconvolution for targeted analysis of data-
18 independent acquisition mass spectrometry proteomics. *Nature Methods*.
19 2018;15(5):371–378.
- 20 [22] Hu A, Lu YY, Bilmes J, Noble WS. Joint Precursor Elution Profile Infer-
21 ence via Regression for Peptide Detection in Data-Independent Acquisition
22 Mass Spectra. *Journal of Proteome Research*. 2019;18(1):86–94.
- 23 [23] Tsou CC, Avtonomov D, Larsen B, Tucholska M, Choi H, Gingras AC,
24 et al. DIA-Umpire: Comprehensive computational framework for data-
25 independent acquisition proteomics. *Nature Methods*. 2015;12(3):258–264.
- 26 [24] Cox J, Mann M. MaxQuant enables high peptide identification rates, indi-
27 vidualized p.p.b.-range mass accuracies and proteome-wide protein quan-
28 tification. *Nature Biotechnology*. 2008;26(12):1367–1372.
- 29 [25] Bertsch A, Gröpl C, Reinert K, Kohlbacher O. OpenMS and TOPP: open
30 source software for LC-MS data analysis. In: *Methods in molecular biology*
31 (Clifton, N.J.). vol. 696. Springer; 2011. p. 353–367.
- 32 [26] Bellew M, Coram M, Fitzgibbon M, Igra M, Randolph T, Wang P, et al.
33 A suite of algorithms for the comprehensive analysis of complex protein
34 mixtures using high-resolution LC-MS. *Bioinformatics*. 2006;22(15):1902–
35 1909.
- 36 [27] Basu S, Davidson I, Wagstaff K. *Constrained clustering: Advances in*
37 *algorithms, theory, and applications*. CRC Press; 2008.

- 1 [28] Sibson R. SLINK: An optimally efficient algorithm for the single-link cluster
2 method. *The Computer Journal*. 1973;16(1):30–34.
- 3 [29] Defays D. An efficient algorithm for a complete link method. *The Computer*
4 *Journal*. 1977;20(4):364–366.
- 5 [30] Ester M, Kriegel HP, Sander J, Xu X. A Density-Based Algorithm for
6 Discovering Clusters in Large Spatial Databases with Noise. In: *Proceed-*
7 *ings of the 2nd International Conference on Knowledge Discovery and Data*
8 *Mining*. vol. 96; 1996. p. 226–231.
- 9 [31] C SR, Michener. A statistical method for evaluating systematic relation-
10 ships. *Univ Kans Sci Bull*. 1958;38:1409–1438. Available from: <http://ci.nii.ac.jp/naid/10011579647/en/>.
- 12 [32] Von Luxburg U, Williamson RC, Guyon I. Clustering: Science or art? In:
13 *Proceedings of ICML Workshop on Unsupervised and Transfer Learning*;
14 2012. p. 65–79.
- 15 [33] Adolfsson A, Ackerman M, Brownstein NC. To cluster, or not to cluster:
16 An analysis of clusterability methods. *Pattern Recognition*. 2019;88:13–26.
- 17 [34] Datta S, Datta S. Comparisons and validation of statistical clustering tech-
18 niques for microarray gene expression data. *Bioinformatics*. 2003;19(4):459–
19 466.
- 20 [35] Shi J, Malik J. Normalized cuts and image segmentation. *IEEE Transac-*
21 *tions on Pattern Analysis and Machine Intelligence*. 2000;22(8):888–905.
- 22 [36] Ng AY, Jordan MI, Weiss Y. On spectral clustering: Analysis and an
23 algorithm. In: *Advances in Neural Information Processing Systems*; 2002.
24 p. 849–856.
- 25 [37] Von Luxburg U. A tutorial on spectral clustering. *Statistics and Comput-*
26 *ing*. 2007;17(4):395–416.
- 27 [38] Borges H, Guibert R, Permiakova O, Burger T. Distinguishing between
28 Spectral Clustering and Cluster Analysis of Mass Spectra. *Journal of Pro-*
29 *teome Research*. 2019;18(1):571–573.
- 30 [39] Cheng Y. Mean shift, mode seeking, and clustering. *IEEE transactions on*
31 *pattern analysis and machine intelligence*. 1995;17(8):790–799.
- 32 [40] Comaniciu D, Meer P. Mean shift: A robust approach toward feature space
33 analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
34 2002;24(5):603–619.
- 35 [41] Schubert E, Rousseeuw PJ. Faster k-Medoids clustering: improving the
36 PAM, CLARA, and CLARANS algorithms. In: *International Conference*
37 *on Similarity Search and Applications*. Springer; 2019. p. 171–187.

- 1 [42] Macqueen J. Some methods for classification and analysis. In: Proceedings
2 of the Fifth Berkeley Symposium on Mathematical Statistics and Probabil-
3 ity, Volume 1: Statistics. vol. 233. Oakland, CA, USA; 1967. p. 281–297.
4 Available from: <http://projecteuclid.org/bsmsp>.
- 5 [43] Lloyd SP. Least Squares Quantization in PCM. IEEE Transactions on
6 Information Theory. 1982;28(2):129–137.
- 7 [44] Jain AK. Data clustering: 50 years beyond K-means. Pattern Recognition
8 Letters. 2010;31(8):651–666.
- 9 [45] Williams CKI. Learning With Kernels: Support Vector Machines, Regu-
10 larization, Optimization, and Beyond. vol. 98. MIT press; 2003.
- 11 [46] Schölkopf B, Smola A, Müller KR. Nonlinear Component Analysis as a
12 Kernel Eigenvalue Problem. Neural Computation. 1998;10(5):1299–1319.
- 13 [47] Henning J, Tostengard A, Smith R. A Peptide-Level Fully Annot-
14 ated Data Set for Quantitative Evaluation of Precursor-Aware Mass
15 Spectrometry Data Processing Algorithms. Journal of Proteome Re-
16 search. 2019;18(1):392–398. Available from: [https://github.com/
17 optimusmoose/ups2GT](https://github.com/optimusmoose/ups2GT).
- 18 [48] Chambers MC, Maclean B, Burke R, Amodei D, Ruderman DL, Neumann
19 S, et al. A cross-platform toolkit for mass spectrometry and proteomics.
20 Nature biotechnology. 2012;30(10):918–920.
- 21 [49] Yu Z, Herman G. On the earth mover’s distance as a histogram similarity
22 metric for image retrieval. IEEE International Conference on Multimedia
23 and Expo, ICME 2005. 2005;2005(2):686–689.
- 24 [50] Courty N, Flamary R, Tuia D. Domain adaptation with regularized opti-
25 mal transport. In: Joint European Conference on Machine Learning and
26 Knowledge Discovery in Databases. Springer; 2014. p. 274–289.
- 27 [51] Majewski S, Ciach MA, Startek M, Niemyska W, Miasojedow B, Gambin A.
28 The wasserstein distance as a dissimilarity measure for mass spectra with
29 application to spectral deconvolution. In: 18th International Workshop
30 on Algorithms in Bioinformatics (WABI 2018). Schloss Dagstuhl-Leibniz-
31 Zentrum fuer Informatik; 2018. .
- 32 [52] Schölkopf B. The kernel trick for distances. In: Advances in Neural Infor-
33 mation Processing Systems; 2001. p. 301–307.
- 34 [53] Wang S, Gittens A, Mahoney MW. Scalable kernel K-means clustering with
35 Nyström approximation: relative-error bounds. The Journal of Machine
36 Learning Research. 2019;20(1):431–479.
- 37 [54] Keriven N, Bourrier A, Gribonval R, Pérez P. Sketching for large-scale
38 learning of mixture models. Information and Inference: A Journal of the
39 IMA. 2018;7(3):447–508.

- 1 [55] Hartigan JA, Wong MA. Algorithm AS 136: A K-Means Clustering Algo-
2 rithm. *Applied Statistics*. 1979;28(1):100.
- 3 [56] Keriven N, Tremblay N, Traonmilin Y, Gribonval R. Compressive K-means.
4 In: *ICASSP, IEEE International Conference on Acoustics, Speech and Sig-
5 nal Processing - Proceedings*. Institute of Electrical and Electronics Engi-
6 neers Inc.; 2017. p. 6369–6373.
- 7 [57] Givens CR, Shortt RM. A class of Wasserstein metrics for probability
8 distributions. *The Michigan Mathematical Journal*. 1984;31(2):231–240.
- 9 [58] Gibbs AL, Su FE. On choosing and bounding probability metrics. *Inter-
10 national Statistical Review*. 2002;70(3):419–435.
- 11 [59] Hofmann T, Schölkopf B, Smola AJ. Kernel methods in machine learning.
12 *Annals of Statistics*. 2008;36(3):1171–1220.
- 13 [60] Berlinet A, Thomas-Agnan C. *Reproducing Kernel Hilbert Spaces in Prob-
14 ability and Statistics*. Springer Science & Business Media; 2004.
- 15 [61] Feragen A, Lauze F, Hauberg S. Geodesic exponential kernels: When
16 curvature and linearity conflict. In: *Proceedings of the IEEE Conference
17 on Computer Vision and Pattern Recognition*; 2015. p. 3032–3042.
- 18 [62] Calandriello D, Rosasco L. Statistical and computational trade-offs in ker-
19 nel K-means. In: *Advances in Neural Information Processing Systems*. vol.
20 2018-Decem; 2018. p. 9357–9367.
- 21 [63] Rahimi A, Recht B. Random features for large-scale kernel machines. In:
22 *Advances in neural information processing systems*; 2008. p. 1177–1184.
- 23 [64] Puckette SE, Rudin W. *Fourier Analysis on Groups..* vol. 72. Wiley Online
24 Library; 1965.
- 25 [65] Arias P, Randall G, Sapiro G. Connecting the out-of-sample and pre-
26 image problems in Kernel methods. In: *Proceedings of the IEEE Computer
27 Society Conference on Computer Vision and Pattern Recognition*. IEEE;
28 2007. p. 1–8.
- 29 [66] Mika S, Schölkopf B, Smola A, Müller KR, Scholz M, Rätsch G. Kernel
30 PCA and de-noising in feature spaces. In: *Advances in Neural Information
31 Processing Systems*; 1999. p. 536–542.
- 32 [67] Prive F, Aschard H, Ziyatdinov A, Blum MGB. Efficient analysis of large-
33 scale genome-wide data with two R packages: Bigstatsr and bigsnpr. *Bioin-
34 formatics*. 2018 aug;34(16):2781–2787. Available from: [https://doi.org/
35 10.1093/bioinformatics/bty185](https://doi.org/10.1093/bioinformatics/bty185).
- 36 [68] Permiakova O, Burger T. Gitlab of CHICKN (Chromatogram Hierarchi-
37 cal Compressive K-means with Nystrom approximation) R package; 2020.
38 Available from: <https://gitlab.com/Olga.Permiakova/chickn>.

1 [69] Permiakova O, Burger T. CRAN repository of CHICKN (Chromatogram
2 Hierarchical Compressive K-means with Nystrom approximation) R pack-
3 age; 2020. Available from: [https://CRAN.R-project.org/package=](https://CRAN.R-project.org/package=chickn)
4 [chickn](https://CRAN.R-project.org/package=chickn).

5 **Figures**

6 **Additional Files**

7 **Additional File 1**

8 **Title:** Ecoli-FMS data matrix.

9 **Description:** Figure depicting the matrix built thanks to the mass spectrum
10 interpolation of Ecoli-FMS data. Each matrix column corresponds to a
11 chromatographic profile for a fixed m/z value. Maximum Intensity for
12 columns and for rows is depicted in bar plots.

13 **Format:** .png file.

14 **Additional File 2**

15 **Title:** Preprocessing details.

16 **Description:** Detailed explanations of Equation 1 (interpolation needs, justi-
17 fication of the method and parameter tuning).

18 **Format:** .pdf file.

19 **Additional File 3**

20 **Title:** Kernel positive (semi-)definiteness.

21 **Description:** Empirical evidences (Gaussian W1 case) and formal demonstra-
22 tion (Laplacian W1 case) of the P(S)D-ness of the proposed kernels.

23 **Format:** .pdf file.

24 **Additional File 4**

25 **Title:** Sketch size influence on the clustering.

26 **Description:** Influence of the sketch size on performances clustering of the
27 Ecoli-DIA dataset, in function of the computational cost and the number
28 of clusters.

29 **Format:** .png file.

1 **Additional File 5**

2 **Title: Kernel hyperparameter stability.**

3 **Description:** Figure showing the stability of the hyperparameter γ of Lapla-
4 cian and Gaussian W1 kernels with respect to the neighborhood maximum
5 size ν .

6 **Format:** .png file.

7 **Additional File 6**

8 **Title: Consensus chromatogram stability**

9 **Description:** A set of 10 figures exemplifying the stability of the pre-image
10 computation through the averaging of a neighborhood of varying size.

11 **Format:** A zipped folder (.zip) containing .png files.

12 **Additional File 7**

13 **Title: Influence of k on the execution time of CHICKN.**

14 **Description:** Figure depicting CHICKN execution time as a function of k ,
15 the number of clusters at each iteration, for both UPS2GT (blue) and
16 Ecoli-FMS (red) datasets.

17 **Format:** .png file.

18 **Additional File 8**

19 **Title: Influence of k_{total} on the execution time of CHICKN.**

20 **Description:** Figure depicting CHICKN execution time as a function of k_{total} ,
21 the maximum number of clusters, for both UPS2GT (blue) and Ecoli-FMS
22 (red) datasets.

23 **Format:** .png file.

24 **Additional File 9**

25 **Title: Performance evaluation for the Laplacian W1 kernel.**

26 **Description:** This figure is the same as Figure 4, yet with $p = 1$ instead of
27 $p = 2$. The performance on the UPS2GT dataset are a bit lower than
28 with the Gaussian W1 kernel (equivalent Rand index, better precision,
29 lower recall), making it unable to compete with Xnet. However, on raw
30 data such as Ecoli-DIA (*i.e.* on data CHICKN should work with), the
31 Laplacian W1 kernel exhibit slightly better DB index than its Gaussian
32 counterpart; however, this is hardly significant, making us conclude that
33 strict performance should not be the criterion to chose the kernel.

1 **Format:** .png file.

2 **Additional File 10**

3 **Title:** Differently charged ions of a same peptide tend to cluster to-
4 **gether.**

5 **Description:** A subset of clusters were manually inspected so as to label as
6 many profiles with the corresponding identified ion. Although this la-
7 belling cannot be exhaustively conducted due to the largely incomplete
8 coverage of MS/MS analysis, it could be established that ions of a same
9 peptide cluster together in many cases.

10 **Format:** .png file.

11 **Additional File 11**

12 **Title:** Cluster size distribution.

13 **Description:** Histograms of the cluster size distribution resulting from the
14 application of CHICKN on each of the three datasets.

15 **Format:** .png file.

16 **Additional File 12**

17 **Title:** Examples of well-formed clusters for the Ecoli-FMS dataset.

18 **Description:** Same figure as Figure 6 with Laplacian W1 kernel.

19 **Format:** .png file.

20 **Additional File 13**

21 **Title:** Examples of multiplexed clusters for the Ecoli-FMS dataset
22 **using CHICKN method.**

23 **Description:** Figure illustrating that dividing multiplexed clusters into several
24 sub-clusters would improve the elution profile interpretation. The real
25 chromatograms and the consensus chromatograms are depicted in gray
26 and in red, respectively.

27 **Format:** .png file.

1 **Additional File 14**

2 **Title:** Differently charged ions of a same peptide tend to cluster to-
3 **gether.**

4 **Description:** Figure similar to Additional File 10. It depicts another subset of
5 CHICKN clusters with chromatographic profiles manually annotated with
6 the corresponding peptide ion. It could be established that ions of a same
7 peptide tend to cluster together.

8 **Format:** .png file.

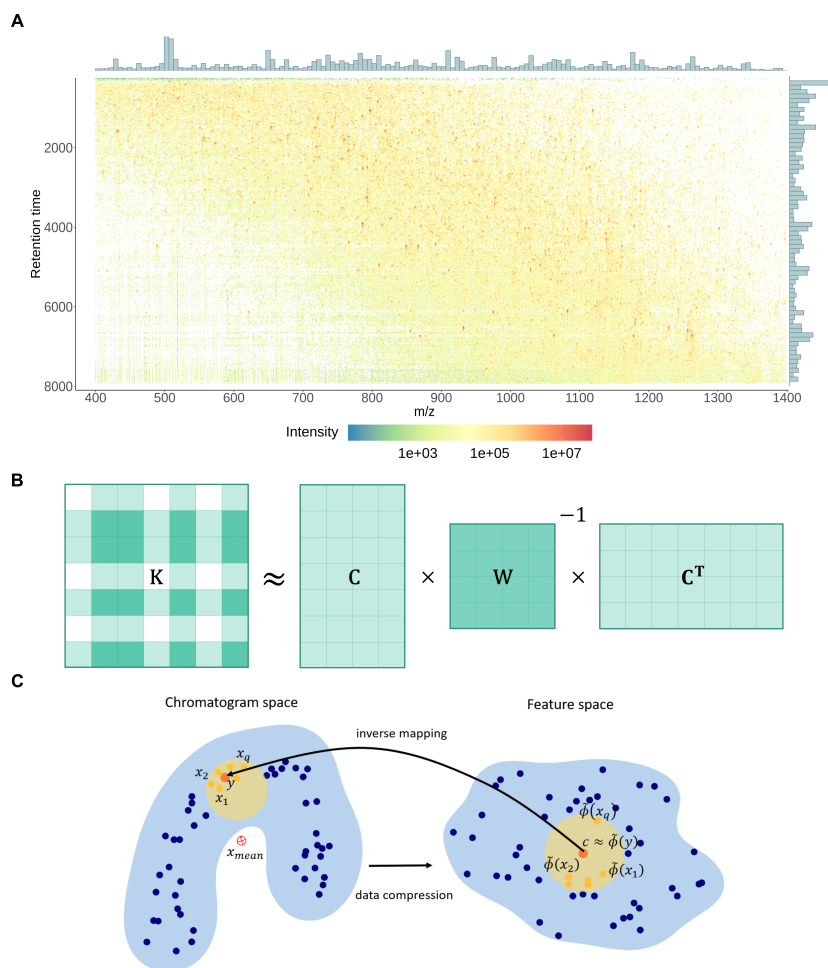


Figure 1: Data matrix, Nyström approximation and pre-image illustrations. (A) Ecoli-DIA data matrix. Each matrix column corresponds to a chromatographic profile for a fixed m/z value. Maximum Intensity for columns and for rows is depicted in bar plots. (B) Nyström kernel approximation. The matrix C represents the similarity between each data point and the random sample. The matrix W corresponds to the pairwise similarity evaluation between selected data points. (C) Pre-image problem. Consensus chromatogram construction amounts to solve a pre-image problem, *i.e.* to map the feature space (right) back to the space of chromatograms (left). Blue points depict the elution profiles (left) and their images in the feature space (right). The red points are the cluster centroid (right) and the corresponding consensus chromatogram (left). The yellow circles represent the cluster centroid and consensus chromatogram neighborhoods. Due to the mapping non-linearity, the mean chromatogram may lie outside the cluster, while the correct consensus chromatogram should belong to it.

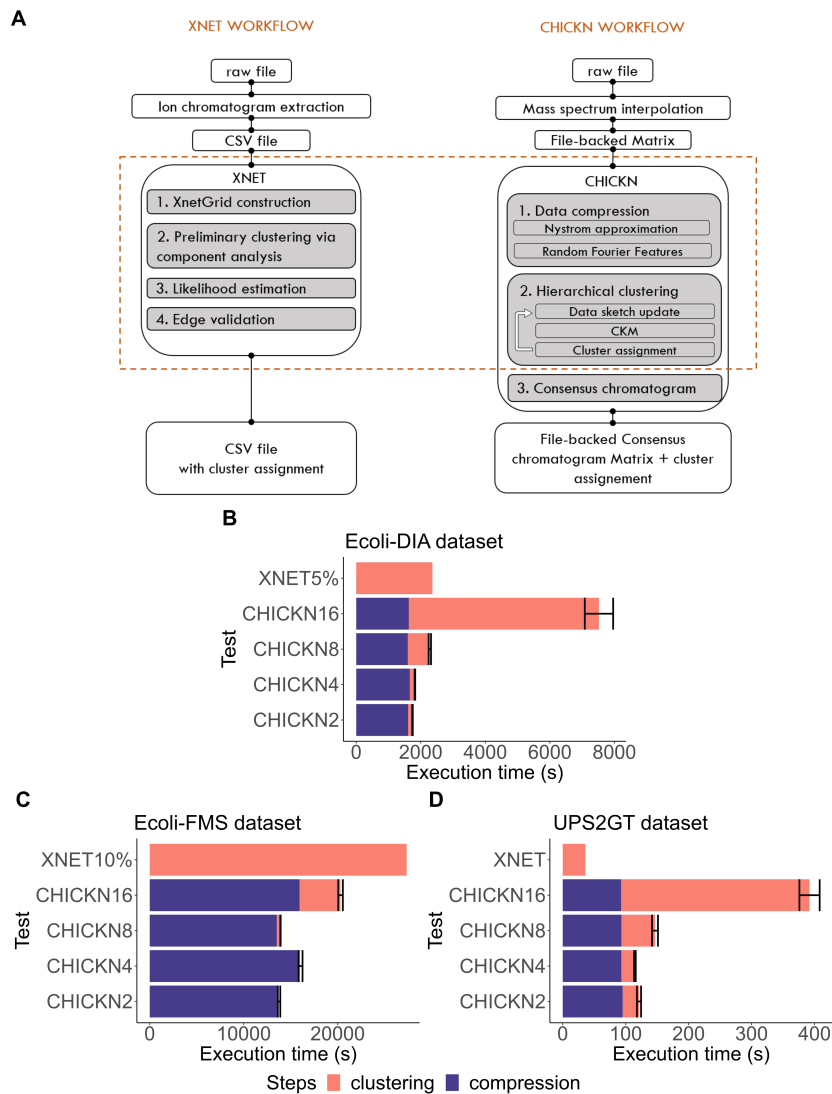


Figure 2: Xnet and CHICKN comparison. (A) The method workflows. To allow for fair comparisons, we have focused on the core algorithms, depicted within the dotted rectangle. (B - D) The execution time comparison for Ecoli and for the UPS2GT datasets. The CHICKN execution time is decomposed into the data compression time (blue) and the clustering time (pink). Note that XNet had to be run on 5% of the Ecoli-DIA dataset and 10% of the Ecoli-FMS dataset only, to avoid "out of memory" issues. The experiments on Ecoli-DIA were performed on a laptop, while other datasets were processed with a multi-core machine.

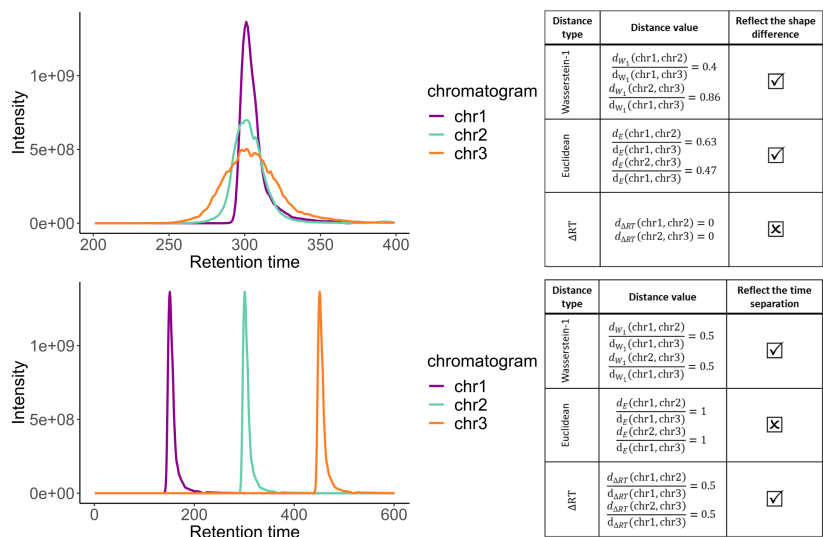


Figure 3: **Distance metrics for chromatographic data analysis.** Comparison of Wasserstein-1, Euclidean and RT difference distances on real chromatographic profiles from the Ecoli-FMS dataset.

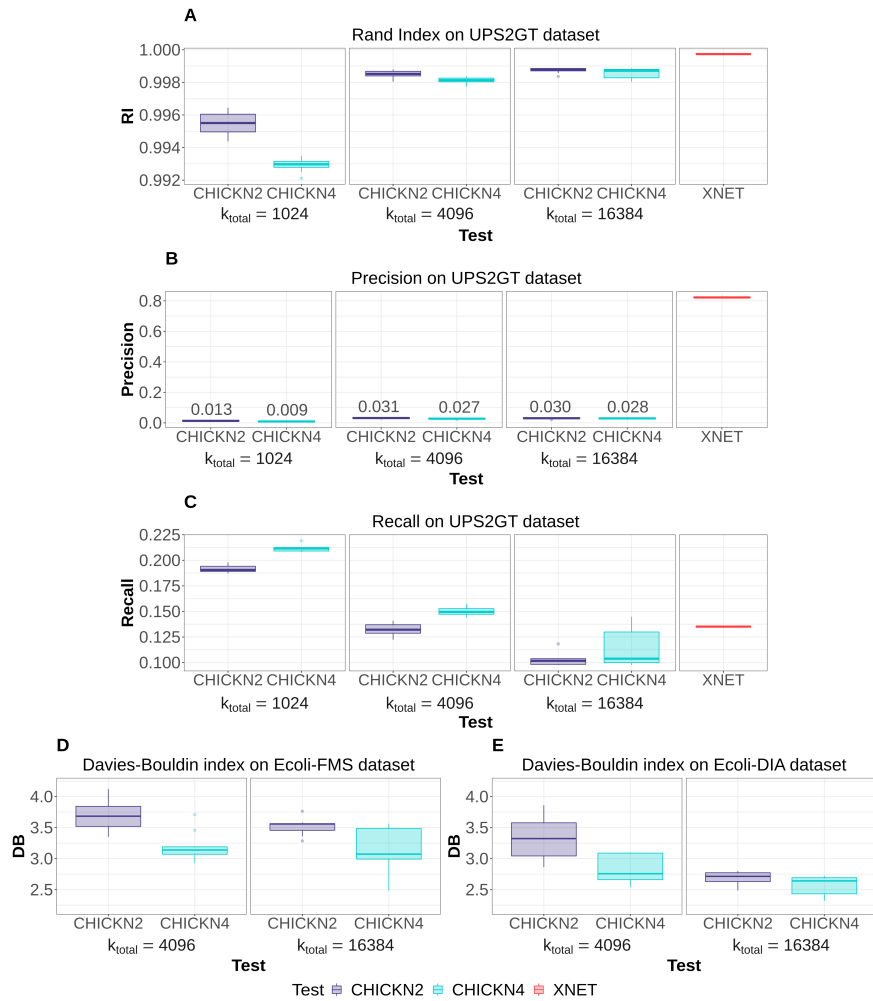


Figure 4: **Statistical result analysis.** (A) Rand index, (B) Precision, (C) Recall and (D-E) DB index depending on the k and k_{total} parameters; CHICKN2 and CHICKN4 tests are depicted in purple and light blue respectively; For the UPS2GT dataset, additional comparisons with Xnet (in red) are provided.

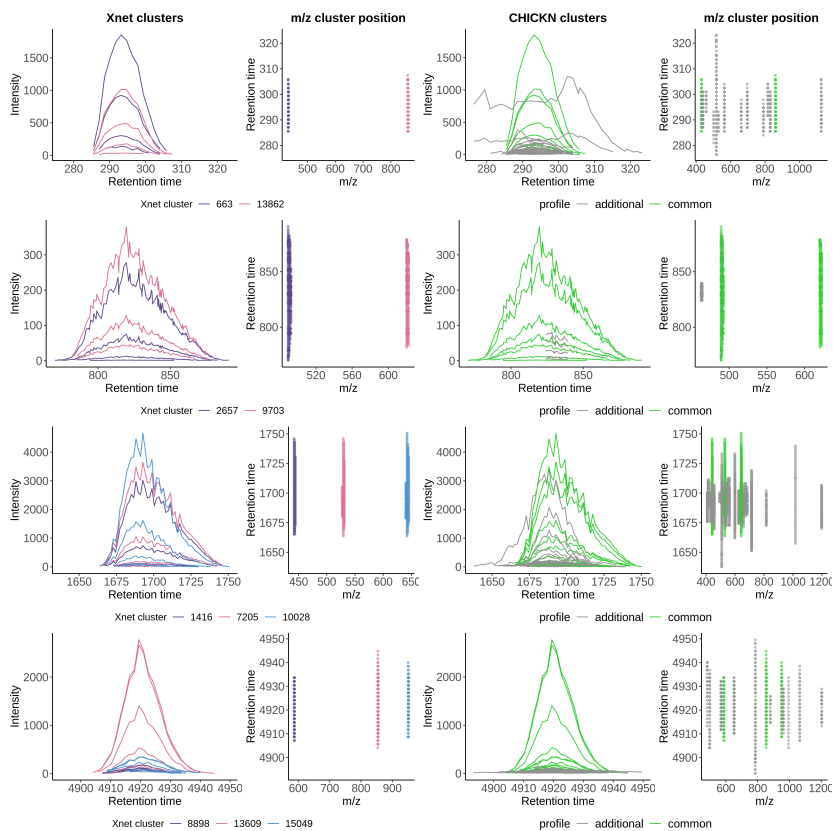


Figure 5: **Xnet and CHICKN clusters for UPS2GT dataset.** Each of the four lines represent a series of chromatograms in the context of their Xnet and CHICKN Cluster. On the plot of the leftmost column, a series of chromatograms with similar shapes are represented in different colors (2 or 3) according to the distinct Xnet clusters they belong to. In the second column, each elution profile is represented with the same color, according to its m/z position, hereby illustrating that Xnet clusters similar signals in different clusters because of a too large m/z difference. The plot of the third column represents the CHICKN cluster which encompasses all the Xnets cluster profiles of the leftmost column (in green), as well as other signals (in gray) falling in the same CHICKN cluster, hereby illustrating CHICK builds meaningful patterns irrespective of the m/z information that is essential to isotopic envelope construction. In the rightmost column, the m/z positions of the signals of the third columns, depicted with the same color code.

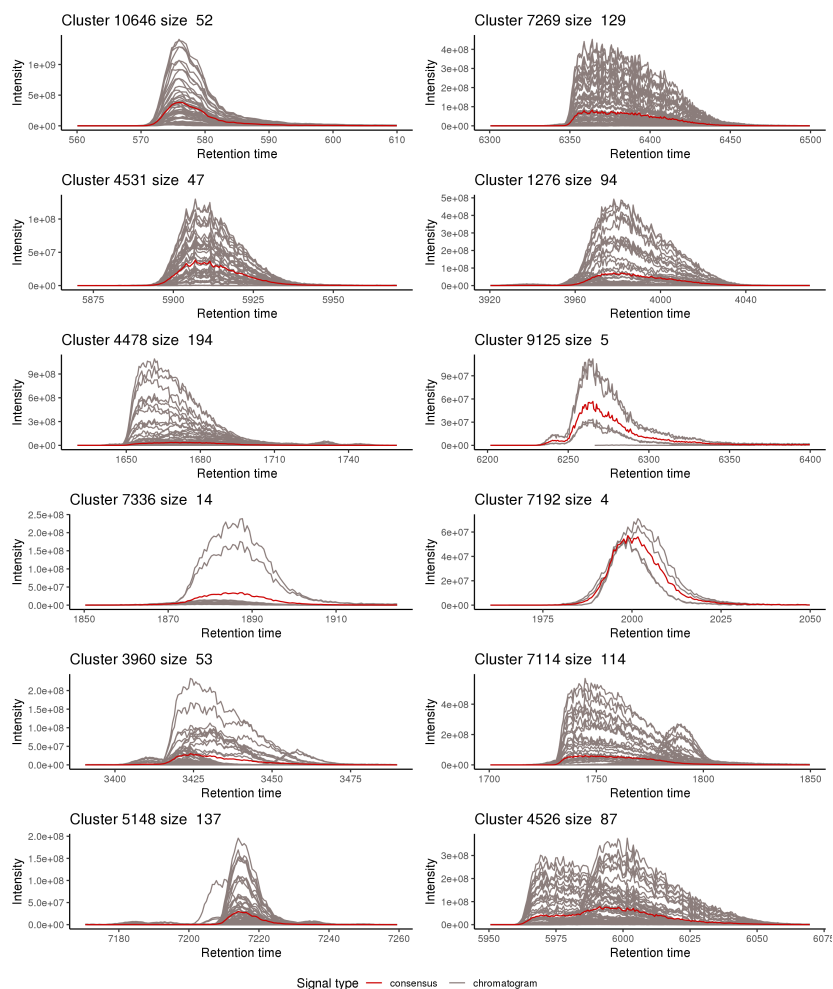


Figure 6: **Examples of well-formed clusters for the Ecoli-FMS dataset.** 12 clusters proposed by CHIKN (represented as time series), where each chromatogram is represented in gray, and where the consensus chromatogram is represented in red. The numbers above each example indicate the cluster ID and the number of chromatograms it encompasses.