

SIMA: Simultaneous Multiple Alignment of LC/MS Peak Lists

Björn Voss[†], Michael Hanselmann[†], Bernhard Y. Renard[‡], Martin S. Lindner, Ullrich Köthe, Marc Kirchner[§] and Fred A. Hamprecht*

Interdisciplinary Center for Scientific Computing, University of Heidelberg, Germany

Associate Editor: Jonathan Wren

ABSTRACT

Motivation: Alignment of multiple liquid chromatography/mass spectrometry (LC/MS) experiments is a necessity today, which arises from the need for biological and technical repeats. Due to limits in sampling frequency and poor reproducibility of retention times, current LC systems suffer from missing observations and non-linear distortions of the retention times across runs. Existing approaches for peak correspondence estimation focus almost exclusively on solving the *pairwise* alignment problem, yielding straightforward but suboptimal results for *multiple* alignment problems.

Results: We propose SIMA, a novel automated procedure for alignment of peak lists from multiple LC/MS runs. SIMA combines hierarchical pairwise correspondence estimation with *simultaneous* alignment and *global* retention time correction. It employs a tailored multidimensional kernel function and a procedure based on maximum likelihood estimation to find the retention time distortion function that best fits the observed data. SIMA does not require a dedicated reference spectrum, is robust with regard to outliers, needs only two intuitive parameters and naturally incorporates incomplete correspondence information. In a comparison with seven alternative methods on four different datasets, we show that SIMA yields competitive and superior performance on real-world data.

Availability: A C++ implementation of the SIMA algorithm is available from <http://hci.iwr.uni-heidelberg.de/MIP/Software>.

Contact: fred.hamprecht@iwr.uni-heidelberg.de

Supplementary information: Supplementary data are available at *Bioinformatics* online.

Received on October 8, 2010; revised on December 10, 2010; accepted on January 24, 2011

1 INTRODUCTION

Recent developments in liquid chromatography/mass spectrometry (LC/MS) have afforded insight into the dynamics of biological systems at unprecedented levels of detail. High-resolution MS-based protein identification and quantitative MS are now established methodologies in fields as diverse as proteomics (Aebersold and Mann, 2003), glycomics (Zaia, 2010), lipidomics (Shevchenko and Simons, 2010) and metabolomics (Dettmer *et al.*, 2007).

*To whom correspondence should be addressed.

[†]The authors wish it to be known that, in their opinion, the first two authors should be regarded as joint First Authors.

[‡]Present address: Department of Computational Medicine. The Institute for Translational Oncology and Immunology, Mainz, Germany.

[§]Present address: Proteomics Center, Departments of Pathology, Children's Hospital Boston and Harvard Medical School, Boston, MA, USA.

Robust alignment of LC/MS experiments: current LC/MS experiments often investigate complex biological systems over a set of different environmental conditions and/or time-courses. The associated data are routinely split into multiple fractions and acquired in technical and biological replicates, yielding tens to hundreds of LC/MS runs. Each of these runs delivers a snapshot of the system of interest and to enable their joint analysis, the common components in different measurements need to be related to each other. In practical LC/MS applications, two major factors complicate the determination of component correspondences across multiple runs: (i) the limited reproducibility attained on LC systems, which gives rise to non-linear distortions of the retention time domain; and (ii) the limited sampling frequency inherent to data-driven MS/MS acquisition as a notorious cause for missing observations. To obtain quantitative estimates or to increase peptide identification rates over a series of experiments, LC/MS data analysis frameworks (Khan *et al.*, 2009; Mueller *et al.*, 2007) rely on accurate mass and retention time alignment to propagate correspondence information between runs, experiments and samples. Accounting for LC distortions is a necessary prerequisite for such cross-experiment inference (America and Cordewener, 2008; Podwojski *et al.*, 2009).

Although numerous pairwise alignment methods have been proposed, the question of *simultaneous* alignment of multiple datasets is still a particularly challenging task: as the number of potential correspondences grows exponentially, false initial multiple correspondence estimates are more likely, and estimation procedures for the associated warping functions need to be robust to potential outliers. Even more, in the light of practical application, multiple alignment methods should naturally cope with incomplete correspondences where peaks are only observed in a subset of runs and no obvious missing value imputation strategy is available.

Types of LC/MS alignment algorithms: published alignment algorithms work on different representations—either peaks extracted from raw measurements or the raw measurements themselves (Clifford *et al.*, 2009; Prakash *et al.*, 2006; Vandenbogaert *et al.*, 2008). This contribution focuses on the alignment of sparse sets of samples, i.e. peaks $p_i = [(m/z)_i, (rt)_i, z_i]$, which lie in a 3D feature space (ion mass-to-charge ratio, ion elution time and ion charge) as proposed in Cox and Mann (2008); Khan *et al.* (2009); Lange *et al.* (2008). Existing approaches can be divided into three categories:

- (1) Alignment based on a static reference, where all observed measurements are aligned to a single reference peak list (Bellew *et al.*, 2006; Lange *et al.*, 2007, 2008; Sturm *et al.*, 2008; Zhang *et al.*, 2005). Because these approaches single out one measured reference run, they perform well if there

is at least one run of exceptional quality. Peaks that are not present in the reference cannot be used for alignment. This can be a substantial drawback, especially in low signal-to-noise ratio (SNR) situations or if suboptimal reproducibility is an issue.

- (2) Complete pairwise correspondence-based alignment, where the pairwise distances between all observed measurements in all runs yield pairwise alignments (Li *et al.*, 2005). Based on these correspondence pairs, global correspondence groups are computed by iteratively linking similar peaks. Although this approach overcomes the single reference problem, it is computationally expensive since it requires the calculation of all similarity measurements between all extracted peaks and performs a retention time correction in each iteration.
- (3) Hierarchical progressive alignment, where a similarity measure based on peak distances determines the merging sequence between different peak lists (Mueller *et al.*, 2007; Prakash *et al.*, 2006). The algorithm starts with an arbitrary (e.g. the most complete) list and subsequently merges the most similar lists until all correspondences are computed. Like in pairwise alignment, the retention times are corrected in each step that is suboptimal (Smith *et al.*, 2006). Few methods exist that work without a similarity measure (Pluskal *et al.*, 2010).

In all categories, existing multiple alignment methods are straightforward extensions of pairwise alignments (Khan *et al.*, 2009; Lange *et al.*, 2007, 2008; Mueller *et al.*, 2007). While some of them use heuristics to deal with peaks that are not present in all available peak lists, i.e. missing correspondences, others completely discard incomplete correspondence information. However, the probability that a peak is visible in *all* peak lists decreases with increasing numbers of runs that have to be aligned. Simultaneous correction of all peak lists is superior but rarely considered (Smith *et al.*, 2006).

Simultaneous multiple LC/MS alignment: we propose SIMA, a novel approach that performs a single global retention time correction based on the multiple correspondence information obtained from all peak lists and naturally deals with missing correspondences. SIMA: (i) uses a pairwise greedy hierarchical strategy to determine all (potentially incomplete) correspondences across D peak lists (without performing a retention time correction in each step; Section 2.1); and then (ii) uses a kernel density estimation type non-parametric approach to simultaneously work on all correspondence groups and derive a D -dimensional *retention time ridge*. Its highest path is found by maximum likelihood estimation and approximates the global *retention time distortion function* that describes retention time shifts across all runs (Section 2.2). Finally, (iii) it uses this function to correct the individual peak lists for retention time shifts and optionally performs a second hierarchical correspondence estimation (Section 2.3). Step (ii) relies on a customized kernel that is inspired by *signal maps* (Prakash *et al.*, 2006) and that has specifically been tailored to make direct use of complete and incomplete correspondence information.

The remainder of this contribution is organized as follows: Section 2 introduces the proposed workflow (see Fig. 1) and its mathematical framework. Section 3 describes the experimental setup and error statistics that is used to judge SIMA performance and Section 4 reports and discusses the outcomes. We end with conclusions in Section 5.

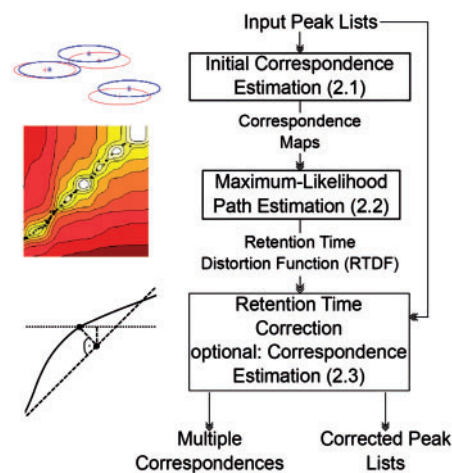


Fig. 1. SIMA workflow: starting from a set of D LC/MS peak lists, SIMA conducts an initial correspondence estimation that yields groups of corresponding peaks (Section 2.1). Based on these groups, the method calculates the retention time distortion function that is most likely to explain the observed retention time differences across the D peak lists (Section 2.2). Given this function, all retention time deviations are corrected, and peak correspondences are reestimated (optional) (Section 2.3).

2 METHODS

2.1 Correspondence estimation

The correspondence estimation is based on: (i) a measure for estimating the distance of peaks from two different peak lists, (ii) an algorithm for establishing peak correspondence pairs based on this measure, (iii) a distance measure for quantifying the dissimilarity of two peak lists based on their peak correspondences and (iv) a hierarchical iteration scheme that successively combines the sparse individual peak lists into a more complete master peak list, while storing all established peak correspondences.

Notation: let $\mathcal{P} = \{P_d\}$, $d = 1, \dots, D$, be the set of the D LC/MS peak lists to be aligned. The d -th peak list P_d comprises $|P_d|$ peaks, $P_d = \{p_{d,1}, \dots, p_{d,|P_d|}\}$, and each peak $p_{d,i} = [(m/z)_{d,i}, (rt)_{d,i}, z_{d,i}]$ is described by its mass over charge position (m/z) , retention time (rt) and charge state z . To simplify notation, we discard the index indicating the membership of a peak to a peak list throughout the remainder of the derivations. Scalars are printed in standard font and vectors in bold.

Distance measure for peak correspondence estimation: the definition of an adequate peak distance measure is fundamental for identifying peak correspondences between different LC/MS runs. We quantify the distance between two peaks p_i and p_j by the diagonal thresholded squared Mahalanobis distance ϕ given by

$$\phi(p_i, p_j) = \Psi((p_i - p_j) * W * (p_i - p_j)') \quad (1)$$

where we define $\Psi(\psi) = \psi$ for $0 \leq \psi \leq 1$, $\Psi(\psi) = \infty$ for $\psi > 1$, and the weight matrix W as $W = \text{diag}^{-1}(T_{(m/z)}^2, T_{(rt)}^2, T_z^2)$. $T_{(m/z)}$ and $T_{(rt)}$ are user-defined threshold parameters for the upper bounds on (m/z) and (rt) shift tolerance. In practice, their choice depends on the measurement precision and is determined by the experimental instrument setup. Choosing very small values for T_z disallows correspondences between peaks with different charge states (default). If reliable charge state information is not available, deviations in charge state may be allowed by using larger values for T_z . Furthermore, W may easily be adapted to also take other features like intensity differences into account (cf. Supplementary Material A). In the 2D $[(m/z), (rt)]$ domain, Equation (1) yields elliptical equidistance lines within the feasible area in which peaks may correspond to each other (cf. Supplementary Material B).

Establishing correspondence groups: given $\phi(\cdot, \cdot)$, the problem of finding correspondences between peaks from two peak lists P_d and P_e can efficiently be solved by an algorithm that is best known for solving the ‘stable marriage’ problem (Gale and Shapley, 1962). Initially designed for graph-matching problems, this method computes an optimal matching between elements from two disjunct sets, such that peak pairs with small distances are preferred. The resulting set \mathcal{F}_{de} contains all peak correspondences of P_d and P_e , that is each peak pair $(i, j) \in \mathcal{F}_{de}$ contains exactly one peak from both P_d and P_e . Note that the two peaks forming a pair may differ in (m/z) , (rt) , and z . Some peaks might not find a partner.

Distance measure for peak list dissimilarity: given \mathcal{F}_{de} , the dissimilarity $\Phi(P_d, P_e)$ of two peak lists is obtained from averaging the finite truncated squared Mahalanobis distances of the assigned peak pairs. Denoting the peaks in pair (i, j) as \mathbf{p}_i and \mathbf{p}_j , we obtain

$$\Phi(P_d, P_e) = \frac{1}{|\mathcal{F}_{de}|} \sum_{(i,j) \in \mathcal{F}_{de}} \phi(\mathbf{p}_i, \mathbf{p}_j). \quad (2)$$

Hierarchical correspondence estimation: rather than relying on one predetermined reference peak list for the alignment, we follow the idea of Mueller *et al.* (2007) and apply a greedy pairwise hierarchical iterative approach: This strategy eliminates the bias toward a single LC/MS run, which occurs when using a reference peak list in the correspondence estimation. Nonetheless, SIMA can also use a single reference peak list to compute multiple peak correspondences. This may be beneficial if one of the peak lists is a priori known to be correct.

We successively combine the peak lists until all individual peak lists have been absorbed in one master peak list. In parallel, we construct *correspondence groups*, i.e. sets of peaks from the individual peak lists that match (see Supplementary Material C for details). Let $\mathcal{P}(t)$ be the set of peak lists that still have to be combined in iteration t . We initialize $\mathcal{P}(0) = \mathcal{P}$, that is with all original peak lists. During the course of the iterations, $\mathcal{P}(t)$ may contain both members of the set of original peak lists and/or representatives for previously combined lists. In each iteration, the two most similar peak lists according to Equation (2) are combined. Assume that in iteration step t these are $P_d \in \mathcal{P}(t)$ and $P_e \in \mathcal{P}(t)$. First, an empty peak list P_{de} is created, and all peaks that are unique to either P_d or P_e are added to it. Then, all peak correspondences $(i, j) \in \mathcal{F}_{de}$ between P_d and P_e are considered, and one representative peak is added for each correspondence pair. Its (rt) and (m/z) values are set to the mean over the respective values of all peaks that in previous iterations have contributed to the two merging peaks. Finally, P_d and P_e are removed from $\mathcal{P}(t)$ and replaced by the combined peak list P_{de} yielding $\mathcal{P}(t+1)$. The correspondence groups are updated accordingly. After $D-1$ steps, all peak lists have been combined, i.e. $|\mathcal{P}(D-1)| = 1$.

We note that the greedy nature of the correspondence estimation allows for an efficient implementation. However, once merged, peaks cannot be split at later iterations, which may suggest that the respective peaks should rather be kept separate. This typically does not pose a practical problem, since by setting the thresholds in W the experimentalist can control the merging behavior of peaks.

Let N be the number of resulting correspondence groups. After iteration $D-1$, the retention times associated with the peaks in the N groups are stored in a *retention time correspondence map* $C \in \mathbb{R}^{N \times D}$. More precisely, element $c_{n,d}$ of C holds the retention time of the peak in P_d that is a member in correspondence group n . If no such peak exists, the respective entry is flagged to indicate a *missing correspondence*. Each row vector $\mathbf{c}_n \in \mathbb{R}^D$, $n = 1, \dots, N$, in C can be interpreted as a *correspondence point* in the D -dimensional *retention time space* (see Figs 4 and 6).

2.2 Maximum likelihood path estimation

Retention time distortion function: define a *master time scale* (MTS) in the retention time space by equidistant sampling of the line of unit slope (that is the angle bisection line for $D=2$). Further, define the *retention time distortion function* (RTDF) as the function that describes the retention time shifts for all peak lists. Its trajectory in retention time space thus explains the observed

correspondence points \mathbf{c}_n , $n = 1, \dots, N$. For a set of perfectly reproducible LC/MS measurements, all \mathbf{c}_n lie on the line of unit slope such that the RTDF is equivalent to the latter. In practice, however, retention time measurements are subject to correlated noise and non-linearly deviate from the ideal case. Given an estimate for the RTDF, the distortion of a peak can be identified by back-projecting its retention time onto the MTS (cf. Fig. 5).

The behavior of the estimate should be in agreement with the fundamental physical properties of LC/MS. We argue that a suitable estimation procedure should: (i) yield a RTDF that features a certain degree of *smoothness* since we do not expect abrupt changes in the elution process, (ii) ensure that the RTDF is *monotonous* such that the elution order is preserved across runs (Kirchner *et al.*, 2007), (iii) be *robust* with respect to measurement errors and naturally deal with outliers that may originate from incorrect matches in the initial correspondence estimation and (iv) be *independent* of the input order of peak lists. We cast the problem of estimating the RTDF into a maximum likelihood (MFL) estimation framework, i.e. we find the RTDF as the function that best explains the correspondence points and at the same time fulfills the above constraints. To this end, we define a customized kernel that incorporates all prior assumptions. Its convolution with the (partially incomplete) observed peak correspondences yields a *retention time ridge*. The path along the highest points of this ‘height profile’, i.e. the maximum likelihood path, is the RTDF that describes the non-linear retention time distortions across the set of peak lists.

Constructing a retention time ridge using a sigmoid kernel: whereas smoothness is guaranteed by employing a smooth kernel, the monotonicity of the ridge in retention times is more difficult to achieve. Figuratively speaking, the kernel (cf. Fig. 2) should induce two preferred areas (lower left, upper right) and two ‘forbidden regions’: when convolving it with all correspondence points in C , the response at a point $\mathbf{x} \in \mathbb{R}^2$ merely depends on the contribution of correspondence points whose retention times are not both lower or both higher than the ones of \mathbf{x} . This encourages the monotonicity of the ridge. Whereas, theoretically speaking, sets of correspondence points can be constructed that lead to paths that violate the monotonicity constraint, all point sets that can be considered a reasonable outcome of a set of LC/MS measurements yield a monotonous result. To make our approach more robust against measurement errors, an adaptive kernel parameter is used that controls the slope and hence the smoothness of the kernel and decreases the influence of outliers. Finally, our method is independent of the input order of the peak lists, since we use an equal kernel profile along all dimensions. By construction, all discussed properties carry over to the higher dimensions.

More formally, let the kernel $K: \mathbf{x} \in \mathbb{R}^D \rightarrow (0, 1)$ be an outer product of sigmoid functions $k(x, \alpha) = 1/(1 + e^{-\alpha x})$ where

$$K(\mathbf{x}) = \prod_{d=1}^D k(x_d, \alpha) + \prod_{d=1}^D k(-x_d, \alpha). \quad (3)$$

Here, $x_d \in \mathbb{R}$ denotes the d -th component of \mathbf{x} , and $\alpha \in \mathbb{R}^+$ is a parameter that controls the influence of the estimated correspondences in the *retention time space* (see Fig. 4c). A higher value for α yields a steeper slope of the sigmoid function $k(\cdot, \cdot)$ and thus increases the local influence of the kernel (see below). For an arbitrary point $\mathbf{x} \in \mathbb{R}^D$, we obtain the cumulative kernel response H with regard to all correspondence points \mathbf{c}_n with the convolution

$$H(\mathbf{x}) = \frac{1}{\Omega} \sum_{n=1}^N \omega(\mathbf{c}_n) K(\mathbf{x} - \mathbf{c}_n) = \quad (4)$$

$$\frac{1}{\Omega} \sum_{n=1}^N \omega(\mathbf{c}_n) \left[\prod_{d=1}^D k(x_d - c_{n,d}, \alpha) + \prod_{d=1}^D k(-x_d + c_{n,d}, \alpha) \right] \quad (5)$$

where $\omega(\mathbf{c}_n)$ is a weighting factor and $\Omega = \sum_{n=1}^N \omega(\mathbf{c}_n)$. To deal with missing correspondences, we replace $k(\cdot, \cdot)$ with the adapted version $\tilde{k}(\cdot, \cdot)$, given by $\tilde{k}(x_d - c_{n,d}, \alpha) = k(x_d - c_{n,d}, \alpha)$ if $c_{n,d} \neq 0$ and 1 otherwise. In both cases, an analytical solution for the derivative $H'(\mathbf{x})$ exists (cf. Supplementary Material D–F). The intuition behind the adaption is as follows: in case of missing

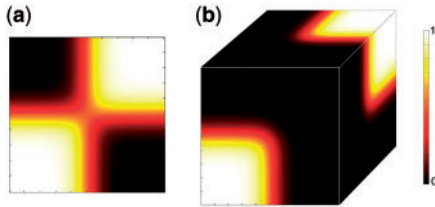


Fig. 2. Plot of the sigmoid kernel defined in Equation (3) for (a) $D=2$ and (b) $D=3$. In both cases, the kernel has two preferred areas along the angle bisection line (light areas) and $(2^D - 2)$ ‘forbidden regions’ (dark areas).

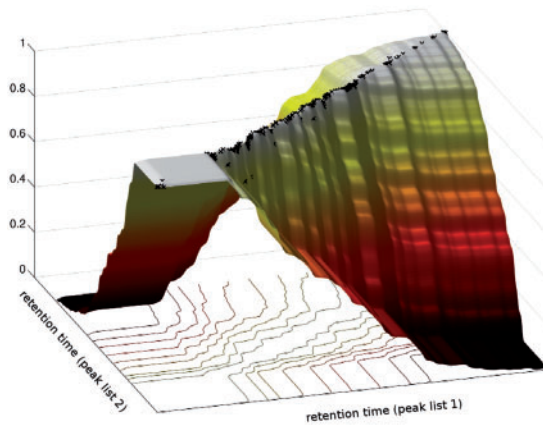


Fig. 3. 3D plot of a retention time ridge formed by convolution of the correspondence points from two peak lists (dots) with the sigmoid kernel (cf. Fig. 2a). Also see movie 3 (Supplementary Material).

correspondences, the correspondence points degenerate to correspondence hyperplanes (cf. Fig. 6). Although incomplete, these correspondences still constrain the RTDF in the orthogonal subspace within the retention time domain. We hence adapt the kernel to be uniform along the missing dimensions. This way, dimensions for which no correspondence information is available are simply ignored, whereas all other information is used whenever available. Here, we use equal weights for all correspondence points, that is $\omega(c_n) = 1 \forall n = 1, \dots, N$. However, different weighting schemes are possible (cf. Supplementary Material G). An exemplary retention time ridge is shown in Figures 3 and 4.

ML-estimation of the retention time distortion function: using the sigmoidal kernel from above, the RTDF can be estimated by finding the points on the highest path through the retention time ridge that correspond to the time points of the MTS. We start with sampling a set of L equidistant points $y_l \in \mathbb{R}^D$, $l = 1, \dots, L$ from the line of unit slope that constitute the MTS, and perform L gradient ascents toward the retention time ridge, starting from each of the y_l . Each gradient ascent is performed on a subspace of \mathbb{R}^D (cf. Fig. 4). These subspaces \mathcal{H}_l are hyperplanes perpendicular to the line of unit slope and contain the y_l , that is, they are described by the normal vector $[1, \dots, 1]$ and support vectors y_l . The gradient ascents yield L points $x_l \in \mathbb{R}^D$ whose piecewise linear interpolation approximates the RTDF. Mathematically, this procedure is similar to a maximum likelihood (ML) estimation where $K(x) \in (0, 1)$ acts as a prior and we determine the x_l by

$$\arg \max_{(x_1, \dots, x_L)} \sum_{l=1}^L H(x_l) \text{ subject to } x_l \in \mathcal{H}_l. \quad (6)$$

Formulas for the normalized gradient directions along which we search for the maxima and for the update of the current estimate of x_l in iteration t are

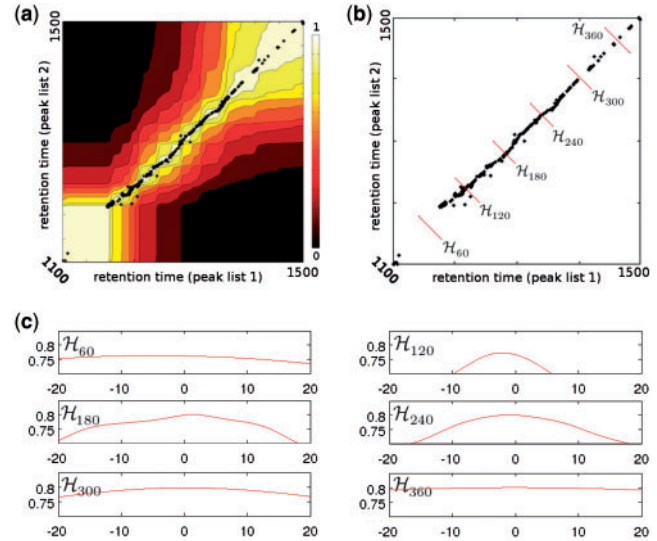


Fig. 4. (a) Contour plot of the retention time ridge (cf. Fig. 3). (b) Correspondence points (dots) and selected hyperplanes \mathcal{H}_l . (c) The intersections of the retention time ridge with these hyperplanes yield height profiles on which the gradient ascents are performed. The profiles of \mathcal{H}_{120} and \mathcal{H}_{240} show distinct bumps. There, the correspondence point density is high, leading to a larger kernel parameter α , i.e. a steeper kernel that emphasizes the influence of local points. The opposite holds, e.g. for \mathcal{H}_{60} .

derived in Supplementary Material D. We propose to use an adaptive step size for the gradient ascents based on the Powell–Wolfe conditions (Powell, 1976) for increased robustness (Supplementary Material F). Note that when performing the gradient ascents, our method never computes the retention time ridge in its entirety but only evaluates $H(x)$ at a few points.

Adaptive kernel bandwidth: naturally, the density of the correspondence points varies throughout retention time space. In such scenarios, the performance of non-parametric kernel methods can be improved by introducing an adaptive kernel bandwidth (Brockmann *et al.*, 1993). Thus, we locally adapt the kernel parameter α as follows: in low-density areas, α is set to low values to achieve a higher robustness and avoid artifacts in the RTDF caused by single observations. In areas of high density, the smoothness of the RTDF is reduced by using larger values for α (cf. Supplementary Material H). This trades off bias and variance and reduces the overall error compared with a non-adaptive scheme.

2.3 Retention time correction

After performing the L gradient ascents and subsequent linear interpolation, we obtain the piecewise linear RTDF which we use for correcting the retention times of the peaks observed in the D peak lists. Assume we want to correct the retention time for peak $p_i = [(m/z)_i, (rt)_i, z_i]$ in peak list P_d . We first find m_1 , the intersection of the RTDF with the hyperplane given by support vector a with $a_d = (rt)_i$ and 0 elsewhere and normal vector $a/\|a\|$. We then identify that point on the line of unit slope that is closest to this intersection point (m_2). The distance of those two points in the d -th dimension constitutes the amount by which the retention time of p_i needs to be corrected. The procedure is repeated for all peaks and all peak lists (see Fig. 5).

Second correspondence estimation: correction of the retention times after the first iteration may give more correspondences. Hence, a second correspondence estimation (cf. Section 2.1) may yield a slightly more complete correspondence map C . However, practical impact is limited due to the overall robust nature of SIMA.

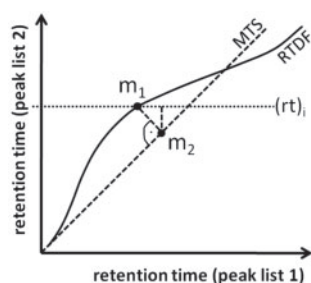


Fig. 5. Retention time correction for $D=2$ where we correct the retention time for peak $p_i = [(m/z)_i, (rt)_i, z_i]$ in peak list P_2 . We obtain m_1 as the intersection of $y = (rt)_i$ with the retention time distortion function (RTDF). By back-projection onto the MTS, we obtain m_2 . The difference in the vertical distance between m_1 and m_2 constitutes the amount by which the retention time of p_i needs to be shifted.

3 EXPERIMENTS

Real-world data: Lange *et al.* (2008) compared a total of six alignment algorithms on four publicly available proteomics ($P1$, $P2$) and metabolomics ($M1$, $M2$) datasets, including msInspect (May *et al.*, 2007), MZmine (Katajamaa and Oresic, 2005; Katajamaa *et al.*, 2006), OpenMS (Lange *et al.*, 2007; Sturm *et al.*, 2008), SpecArray (Li *et al.*, 2005), XAlign (Zhang *et al.*, 2005) and XCMS (Smith *et al.*, 2006). In addition, Pluskal *et al.* recently proposed MZmine 2 with RANSAC aligner (Pluskal *et al.*, 2010). We obtained the proteomics datasets used by Lange *et al.* (2008) from the Open Proteomics Database (OPD) (Prince *et al.*, 2004). Dataset $P1$ originates from an *Escherichia coli* sample and contains two LC/MS runs of six fractions. Dataset $P2$ represents three LC/MS runs of five fractions of different cell states of *Mycobacterium smegmatis*. The datasets were analyzed by LC/MS/MS on an ESI ion trap mass spectrometer (ThermoFinnigan Dexta XP Plus), exported in centroid mode and preprocessed using TOPP tools (Kohlbacher *et al.*, 2007) resulting in a peak list of (m/z) and (rt) positions, which served as input for all alignment procedures. Each run of a fraction contains between 400 and 5800 peaks. Lange *et al.* (2008) optimized parameters for all approaches on the first fraction of each dataset and generated a partial ground truth by linking MS/MS search results from SEQUEST to the LC/MS spectra. A detailed description of all steps and the parameterization of the algorithms is given in Supplementary Material I and (Lange *et al.*, 2008).

For the metabolomics samples, Lange *et al.* (2008) analyzed *Arabidopsis thaliana* leaf tissue using an API QSTAR Pulsar i (Applied Biosystems/MDS Sciex) for the $M1$ dataset and a MicrOTOF-Q (Bruker Daltonics) for the $M2$ dataset resulting in 44 and 24 LC/MS spectra, respectively. Peaks were identified using XCMS (Smith *et al.*, 2006) resulting in 4000 to 17 600 data points per LC/MS spectrum. They generated ground truth by identifying highly confident peak groups, which were reproducible over at least four runs and did not only have the same retention time, but also showed high correlation in their chromatographic peak shapes. Parameters were optimized on the complete datasets for all algorithms since no separate fractions were available. Even though SIMA can operate without a predetermined alignment order, it was run in the same starwise manner (i.e. using one predefined reference against which all remaining runs were aligned) that was used for the other algorithms and the ground truth generation in order to enable

a fair comparison (Otherwise, SIMA might benefit from *not* using a potentially incomplete reference). Again we refer to Lange *et al.* (2008) and the Supplementary Materials for a detailed description and the parameterization of the algorithms.

Performance measures: to measure the performance of an approach, we compute its precision (PR) and recall (RE). Precision is the fraction of correctly aligned peaks among all peaks aligned by one approach, $PR = \frac{\# \text{correctly aligned peaks}}{\# \text{aligned peaks}}$, whereas recall corresponds to the fraction of correctly aligned peaks by one approach among all correct peaks according to the ground truth, $RE = \frac{\# \text{correctly aligned peaks}}{\# \text{correct peaks}}$. To simplify comparison, we use the F-measure $F = \frac{2 \cdot PR \cdot RE}{PR + RE}$, which summarizes precision and recall value by computing their harmonic mean (Gay *et al.*, 2002).

4 RESULTS AND DISCUSSION

SIMA yields competitive or superior results: in Table 1, the results of the comparison are detailed. With regard to recall, SIMA and MZmine 2 (RANSAC) tie at the best performance with an average recall of 0.90. While MZmine 2 (RANSAC) and OpenMS feature better recall values for the $P1$ and $P2$ datasets, SIMA shows better recall performance on the $M1$ and $M2$ data. With regard to precision, i.e. the likelihood of results being correct, SIMA shows the best performance in all datasets with an average precision value of 0.81. This is also reflected in the F-measure, which combines recall and precision. Here, SIMA shows the best overall performance with an average value of 0.85. SIMA always is among the best two methods with respect to the F-measure ($P1$, $P2$) or even performs best ($M1$, $M2$). A complete list of the results on all fractions of all datasets is given in Supplementary Material J.

It is important to note that the comparison favors the algorithms that require a reference peak list: the ground truth generated by Lange *et al.* (2008) is based on the same reference run as used for the alignment. For independently generated ground truth, results for any reference spectrum-based approach are bound to deteriorate since not all peaks present in the ground truth necessarily need to be in the reference peak list. The performance measurements of hierarchical approaches such as SIMA are not affected by this kind of ground truth generation.

SIMA is especially powerful when aligning numerous spectra: The strength of the SIMA approach is particularly visible on the metabolomics ($M1$ and $M2$) datasets. Here, it outperforms the other methods with regard to precision as well as recall. The metabolomics datasets contain more LC/MS spectra (44 and 24, respectively) than the proteomics datasets (2 and 3, respectively) and, thus, also show significantly more missing correspondences. Further, visual inspection confirms that these datasets are less perturbed by noise and show a more characteristic structure, which benefits more from the non-linear fitting of SIMA than the proteomics set, for which linear methods already show good results.

Exploiting incomplete correspondence information is feasible: visual examples of SIMA alignment results are given in Figure 6 as well as in Supplementary Material K and movies 1 and 2. Figure 6 and movie 1 show an (m/z) 500–800 subrange of the first three peak lists in the $M1$ dataset, in which 8 correspondence groups are complete, and incomplete information is available for an additional 14 groups (omitting single entry correspondence groups for visual clarity). The exploitation of partial correspondence groups yields

Table 1. Comparison of the results of seven current alignment approaches with SIMA based on the datasets of the comparative studies by Lange *et al.* (2008) and Pluskal *et al.* (2010)

Data	Measure	ms-Inspect	MZ-mine	Open-MS	Spec-Array	X-Align	XCMS	MZ-mine 2 (RANSAC)	SIMA
P1	RE	0.66	0.85	0.93	0.70	0.88	0.81	0.94	0.92
	PR	0.50	0.89	0.93	0.70	0.88	0.80	0.94	0.94
	F	0.57	0.87	0.93	0.70	0.88	0.80	0.94	0.93
P2	RE	0.58	0.77	0.83	0.50	0.73	0.70	0.75	0.76
	PR	0.26	0.66	0.72	0.35	0.63	0.59	0.68	0.72
	F	0.36	0.71	0.77	0.41	0.67	0.64	0.71	0.74
M1	RE	0.27	0.89	0.87	–	0.88	0.94	0.91	0.92
	PR	0.46	0.74	0.69	–	0.70	0.70	0.74	0.75
	F	0.34	0.81	0.77	–	0.78	0.80	0.82	0.83
M2	RE	0.23	0.98	0.93	–	0.93	0.98	0.98	0.99
	PR	0.47	0.84	0.79	–	0.79	0.78	0.83	0.84
	F	0.31	0.90	0.85	–	0.85	0.87	0.90	0.91
All	RE	0.43	0.87	0.89	–	0.85	0.86	0.90	0.90
	PR	0.42	0.78	0.78	–	0.75	0.72	0.80	0.81
	F	0.39	0.82	0.83	–	0.80	0.78	0.84	0.85

Recall (RE), precision (PR) and the *F*-measure (F) are reported as an average over various runs on two proteomics (P1, P2) and two metabolomics (M1, M2) datasets as well as an overall average of all datasets (All). Bold print highlights the overall best values for each dataset. MZ-mine 2 (RANSAC) and SIMA show the best overall recall, while SIMA features the highest values for precision and the *F*-measure.

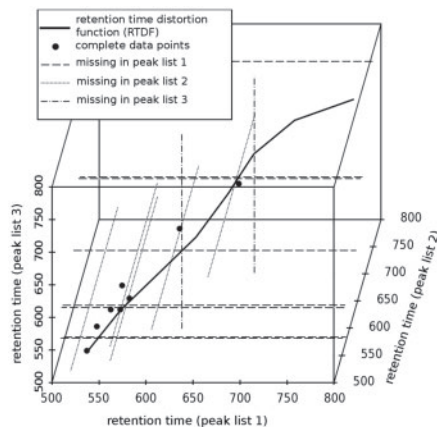


Fig. 6. Visualization of the alignment approach on the peak lists of the first three spectra of the *M1* dataset for the range of 500–800 s in retention time: only 8 peaks (dots) in this range could be matched in all peak lists, whereas 14 peaks were additionally available in two of the three peak lists only. Since these peaks contain information for two dimensions, but are non-informative for the third, they are displayed by straight lines parallel to the coordinate axis of the missing information. The computed *retention time distortion function* (red) still benefits from these straight lines as they help pinpointing its optimal path through this area of few observations.

valuable constraints for guiding the RTDF curve through retention time space and provides robustness in cases where single complete retention time observations show extreme values. The latter is especially prevalent with the obvious outliers present in the mass range (*m/z*) 800–1100, as illustrated in Supplementary Material K and movie 2. SIMA uses all information available from the data by including missing correspondences and can thus base estimates on larger effective numbers of observations compared with other approaches (also cf. Supplementary Material K and L).

Hierarchical correspondence estimation renders distinguished reference spectra obsolete: SIMA eliminates the problem of selecting a distinguished reference spectrum or peak list, respectively. In practical applications, obvious reference candidates are neither easily obtained nor guaranteed to exist. Consequently, SIMA-based alignment is not subject to a reference bias and independent of the peak list processing order.

SIMA is robust with regard to parameter settings: considering that for the proteomics datasets the first fraction was used for parameter optimization, it is interesting to observe that SIMA is not performing as well as, e.g. OpenMS on these fractions. Still, SIMA shows superior performance on the remaining fractions, for which the parameters identified on the first sections were used (cf. Supplementary Material J). This indicates that SIMA is not overly dependent on parameter settings since it not only does benefit from the overfitting on the first fractions where parameters were adjusted to give optimal results, but also shows high-quality results on datasets where parameter optimization was not performed. This can at least partially be explained by the fact that (given reliable charge state information) SIMA only requires two parameters in the matrix *W* for correspondence estimation, which in our experiments have proven robust to changes. Choosing a larger region for the correspondence estimation results in additional random data points for the kernel regression. However, as long as those additional points are unstructured, they do not bias the estimate for the RTDF, since our approach is robust to outliers. Choosing a smaller region results in fewer data points and additional missing correspondences, which can be handled as long as not all data points of a region are removed.

5 CONCLUSION

We introduced SIMA, a novel approach for the simultaneous alignment of multiple LC/MS peak lists. SIMA is specifically tailored to the problems arising from large-scale experiments where

only few peaks are consistently present in all runs. Thus, in contrast to many competing algorithms, SIMA can naturally handle missing correspondences. In addition, it does not rely on a single, error-free reference run as basis for an alignment, but weights the inherent measurement errors of each run against each other.

SIMA requires only very limited user interaction, since it is robust with respect to its two parameters, the thresholds for the tolerated retention time and m/z difference between two peaks. Moreover, these parameters can typically directly be inferred from the expected measurement error in the experiment.

An experimental comparison on real-world proteomics and metabolomics data to seven state-of-the-art approaches demonstrates excellent performance of our method. While matching the recall of MZmine 2 (RANSAC), the best performing method from previous comparisons (Lange *et al.*, 2008; Pluskal *et al.*, 2010), it delivers the best precision and overall F -score values.

Conceptually, SIMA is not limited to the alignment of LC/MS data: by redefinition of the thresholded squared Mahalanobis distance function, it can easily be adapted to any time series with discrete events (cf. Supplementary Material A). SIMA is freely available from <http://hci.iwr.uni-heidelberg.de/MIP/Software>.

ACKNOWLEDGEMENTS

We like to thank Bernhard Kausler, Anna Kreshuk, Sebastian Boppel, and Xinghua Lou (all University of Heidelberg) for fruitful discussions, as well as both of our reviewers for constructive criticism and comments.

Funding: We gratefully acknowledge financial support by the DFG under grant no. (HA4364/2-1 to M.H., B.Y.R.); the Robert Bosch GmbH (to F.A.H.); DFG GRK-1114 (to B.V.).

Conflict of Interest: none declared.

REFERENCES

- Aebersold, R. and Mann, M. (2003) Mass spectrometry-based proteomics. *Nature*, **422**, 198–207.
- America, A.H.P. and Cordewener, J.H.G. (2008) Comparative LC-MS: a landscape of peaks and valleys. *Proteomics*, **8**, 731–749.
- Bellew, M. *et al.* (2006) A suite of algorithms for the comprehensive analysis of complex protein mixtures using high-resolution LC-MS. *Bioinformatics*, **22**, 1902–1909.
- Brockmann, M. *et al.* (1993) Locally adaptive bandwidth choice for kernel regression estimators. *J. Am. Stat. Assoc.*, **88**, 1302.
- Clifford, D. *et al.* (2009) Alignment using variable penalty dynamic time warping. *Anal. Chem.*, **81**, 1000–1007.
- Cox, J. and Mann, M. (2008) MaxQuant enables high peptide identification rates, individualized p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nat. Biotechnol.*, **26**, 1367–1372.
- Dettmer, K. *et al.* (2007) Mass spectrometry-based metabolomics. *Mass Spectrom. Rev.*, **26**, 51–78.
- Gale, D. and Shapley, L. (1962) College admissions and the stability of marriage. *Am. Math. Mon.*, **69**, 15, 9.
- Gay, S. *et al.* (2002) Peptide mass fingerprinting peak intensity prediction: extracting knowledge from spectra. *Proteomics*, **2**, 1374–1391.
- Katajamaa, M. and Oresic, M. (2005) Processing methods for differential analysis of LC/MS profile data. *BMC Bioinformatics*, **6**, 179.
- Katajamaa, M. *et al.* (2006) MZmine: toolbox for processing and visualization of mass spectrometry based molecular profile data. *Bioinformatics*, **22**, 634–636.
- Khan, Z. *et al.* (2009) Protein quantification across hundreds of experimental conditions. *Proc. Natl Assoc. Sci USA*, **106**, 15544–15548.
- Kirchner, M. *et al.* (2007) amsrpm: robust point matching for retention time alignment of LC/MS data with R. **18**, 12.
- Kohlbacher, O. *et al.* (2007) TOPP - the OpenMS proteomics pipeline. *Bioinformatics*, **23**, e191–e197.
- Lange, E. *et al.* (2007) A geometric approach for the alignment of liquid chromatography mass spectrometry data. *Bioinformatics*, **23**, 273–281.
- Lange, E. *et al.* (2008) Critical assessment of alignment procedures for LC-MS proteomics and metabolomics measurements. *BMC Bioinformatics*, **9**, 375.
- Li, X. *et al.* (2005) A software suite for the generation and comparison of peptide arrays from sets of data collected by liquid chromatography-mass spectrometry. *Mol. Cell Proteomics*, **4**, 1328–1340.
- May, D. *et al.* (2007) A platform for accurate mass and time analyses of mass spectrometry data. *J. Prot. Res.*, **6**, 2685–2694.
- Mueller, L.N. *et al.* (2007) SuperHirn - a novel tool for high resolution LC-MS-based peptide/protein profiling. *Proteomics*, **7**, 3470–3480.
- Pluskal, T. *et al.* (2010) MZmine 2: modular framework for processing, visualizing, and analyzing mass spectrometry-based molecular profile data. *BMC Bioinformatics*, **11**, 395.
- Podwojski, K. *et al.* (2009) Retention time alignment algorithms for LC/MS data must consider non-linear shifts. *Bioinformatics*, **25**, 758–764.
- Powell, M.J.D. (1976) Some global convergence properties of a variable metric algorithm for minimization without exact line searches. In Cottle, R.W. and Lemke, C.E. (eds) *Nonlinear Programming, SIAM-AMS Proceedings Volume IX*, SIAM, Philadelphia, PA.
- Prakash, A. *et al.* (2006) Signal maps for mass spectrometry-based comparative proteomics. *Mol. Cell Proteomics*, **5**, 423–432.
- Prince, J.T. *et al.* (2004) The need for a public proteomics repository. *Nat. Biotechnol.*, **22**, 471–472.
- Shevchenko, A. and Simons, K. (2010) Lipidomics: coming to grips with lipid diversity. *Nat. Rev. Mol. Cell Biol.*, **11**, 593–598.
- Smith, C.A. *et al.* (2006) XCMS: processing mass spectrometry data for metabolite profiling using nonlinear peak alignment, matching, and identification. *Anal. Chem.*, **78**, 779–787.
- Sturm, M. *et al.* (2008) OpenMS - an open-source software framework for mass spectrometry. *BMC Bioinformatics*, **9**, 163.
- Vandenbogaert, M. *et al.* (2008) Alignment of LC-MS images, with applications to biomarker discovery and protein identification. *Proteomics*, **8**, 650–672.
- Zaia, J. (2010) Mass spectrometry and glycomics. *OMICS*, **14**, 401–418.
- Zhang, X. *et al.* (2005) Data pre-processing in liquid chromatography-mass spectrometry-based proteomics. *Bioinformatics*, **21**, 4054–4059.