

RESEARCH

Open Access



# Improved de novo peptide sequencing using LC retention time information

Yves Frank, Tomas Hruz, Thomas Tschager\*  and Valentin Venzin

## Abstract

**Background:** Liquid chromatography combined with tandem mass spectrometry is an important tool in proteomics for peptide identification. Liquid chromatography temporally separates the peptides in a sample. The peptides that elute one after another are analyzed via tandem mass spectrometry by measuring the mass-to-charge ratio of a peptide and its fragments. De novo peptide sequencing is the problem of reconstructing the amino acid sequences of a peptide from this measurement data. Past de novo sequencing algorithms solely consider the mass spectrum of the fragments for reconstructing a sequence.

**Results:** We propose to additionally exploit the information obtained from liquid chromatography. We study the problem of computing a sequence that is not only in accordance with the experimental mass spectrum, but also with the chromatographic retention time. We consider three models for predicting the retention time and develop algorithms for de novo sequencing for each model.

**Conclusions:** Based on an evaluation for two prediction models on experimental data from synthesized peptides we conclude that the identification rates are improved by exploiting the chromatographic information. In our evaluation, we compare our algorithms using the retention time information with algorithms using the same scoring model, but not the retention time.

**Keywords:** Computational proteomics, Peptide identification, De novo peptide sequencing, Liquid chromatography, Mass spectrometry

## Background

The amino acid sequences of peptides in a sample can be analyzed by liquid chromatography coupled with tandem mass spectrometry (LC–MS/MS, [1]). First, the peptides are separated temporally by liquid chromatography. Then, the mass spectrometer measures the mass-to-charge ratio of a peptide and fragments multiple copies of it at random positions. Finally, the mass spectrometer measures the mass-to-charge ratio of the resulting fragments. Peptide sequencing [2, 3] is the problem of reconstructing the amino acid sequence of the peptide. When analyzing unknown peptides the otherwise very successful database search approach is not applicable. We focus on de novo sequencing, that is the reconstruction of the

whole amino acid sequence from scratch without the help of a database of known sequences.

Several algorithms for de novo sequencing [4–8] consider the differences of the peptide's fragment masses to reconstruct the peptide's sequence. Various scoring functions have been proposed that try to exploit as much information as possible from the mass spectrum of the fragments to find a sequence that explains the observed spectrum in the best possible way. However, the information obtained from the chromatographic separation in the first step of the LC–MS/MS experiment is not considered by these scoring functions.

In liquid chromatography, the peptides in a sample have to pass through a column. The time a peptide needs to traverse the column is called *retention time* and depends on certain chemical properties of the peptide. This process results in the temporal separation of the peptides in a sample. Predicting the retention time of a peptide from its amino acid sequence is a challenging

\*Correspondence: [tschager@inf.ethz.ch](mailto:tschager@inf.ethz.ch)  
Department of Computer Science, ETH Zurich, Universitätsstrasse 6,  
8092 Zürich, Switzerland



task [9, 10]. Several studies use retention time prediction models for peptide sequencing as a filtering step after a database search to increase the confidence of identification and to identify false positive identifications [11, 12].

However, to the best of our knowledge, the retention time information has not been considered by de novo peptide sequencing algorithms. The retention time can be useful, because it contains information about parts of a sequence that cannot be resolved by mass spectrometry (e.g. amino acids and fragments with equal masses, but different retention times). Moreover, it is available without additional experimental effort. However, simply filtering the candidate sequences of standard de novo sequencing algorithms by their predicted retention time is not an option, as this approach requires to compute all possible candidate sequences in the worst case to find an optimal solution. We formulate and study a de novo sequencing problem that integrates the retention time as an additional constraint and does not require filtering many candidates. We are interested in a sequence that both matches the experimental spectrum and the measured retention time. We consider three additive retention time prediction models and develop algorithms for each model.

In this study,<sup>1</sup> we do not aim for a replacement for available de novo sequencing tools, but rather explore ways of exploiting the retention time information in de novo sequencing algorithms. In the experimental evaluation, we are primarily interested in the impact of using the retention time information. We compare the identification rates of proposed algorithms for two prediction models with the identification rates of DeNovo $\Delta$  [14], an algorithm that uses the same symmetric difference scoring model, but no retention time information. The symmetric difference scoring model already shows improved identification rates compared to the prevalent shared peak count scoring model [5] and this is further improved considering the retention time. We intentionally consider a very basic scoring function to clearly expose the impact of exploiting the retention time information. We evaluate the performance of our algorithms on experimental data of synthesized peptides from the SWATH MS gold standard (SGS, [15]) dataset. For the third prediction model, we present some exemplary results and discuss factors that can limit its applicability. A proof-of-concept implementation of our algorithms is available at Github and can be integrated in the OpenMS framework [16].

Considering the retention time information comes at the cost of higher computational effort and requires additional parameters for retention time prediction. These parameters depend on the chosen standard operating

protocol (SOP) chosen for the experiment and on the LC column of the experiment. Estimating these parameters requires suitable datasets, unless they are available in the literature. Yet, we believe that it is useful to exploit retention time information for peptide identification and to further study the integration of retention time information in algorithms for de novo peptide sequencing.

## Problem definition

### Remarks on model simplifications

To focus on algorithmic aspects of the problem, we simplify several characteristics of the experimental data in our modeling of the de novo peptide sequencing problem. First, the peptide molecule contains an H<sub>2</sub>O molecule in addition to the amino acid molecules. Therefore, the peptide mass has an offset of 18 Da compared to the sum of the amino acid masses. To simplify the description of the algorithms, we do not consider this offset and assume that the mass of a peptide is the sum of the masses of its amino acids. Similarly, we do not consider the fragment mass offsets of different ion types in the description. However, we do consider both offsets in the implementation of our algorithms using techniques described in [14].

Moreover, the mass spectrometer measures mass-to-charge ratios, while our model requires masses as input. Charge state deconvolution [1] is required as a preparatory step to convert mass-to-charge ratios to masses if fragments with a higher charge state should be considered.

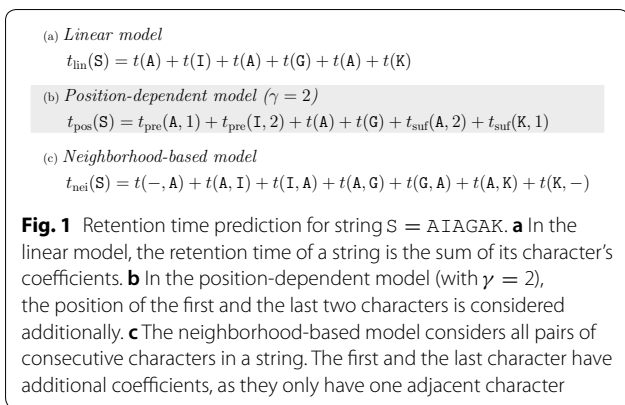
While we do not explicitly model post-translational modifications, our model can consider both fixed and variable modifications. Fixed modifications can be considered by altering the amino acid masses accordingly. Variable modifications are supported by adding new characters to the alphabet of amino acids.

Finally, we consider integer values for the fragment masses and retention times in the description of the algorithm and ignore the mass accuracy of the mass spectrometer. We account for the mass accuracy of the instrument by multiplying the masses by an appropriate factor before rounding to integers. Additionally, in the implementation of our algorithm we consider masses to be equal if they differ at most by a predefined error tolerance (0.02 Da in our experiments).

### Notation

We model an amino acid by a character of an alphabet  $\Sigma$  and a peptide by a string  $S = a_1 \dots a_n$  over  $\Sigma$ . The empty string is denoted by  $S_\emptyset$ . Every character  $a \in \Sigma$  has a mass  $m(a) \in \mathbb{N}$ . The mass of a string  $S = a_1 \dots a_n$  is the sum of its character's masses  $m(S) := \sum_{i=1}^n m(a_i)$ . The empty string  $S_\emptyset$  has mass 0. A substring of  $S$  is denoted

<sup>1</sup> A preliminary version has been presented at WABI 2017 [13].



by  $S_{i,j} = a_i \dots a_j$  for  $1 \leq i \leq j \leq n$ . The prefix set  $\text{Pre}(S)$  contains all prefixes of  $S$  including the empty string, i.e.  $\text{Pre}(S) := \cup_{i=1}^n S_{1,i} \cup \{S_\emptyset\}$ . The *theoretical spectrum* of  $S$  is the union of all its prefix and suffix masses  $\text{TS}(S) := \{m(T), m(S) - m(T) \mid T \in \text{Pre}(S)\}$ . Note that for every prefix  $T \in \text{Pre}(S)$  the string  $S$  has a complementary suffix of mass  $m(S) - m(T)$ . A mass  $m$  is *explained* by  $S$  if  $m \in \text{TS}(S)$ .

### Retention time prediction models

We define three simple models for predicting the retention time of a string  $S = a_1 \dots a_n$  (Fig. 1). The first model is a simple additive model with one retention time coefficient for each character in  $\Sigma$ . The model only considers the character frequencies of a string and has been proposed by [17]. It served as starting point for the development of more evolved prediction models [18, 19].

The other two models consider additional factors that affect the retention time of a peptide. Besides the character frequency, the position of the characters in the string is especially important for the first and the last few positions in the string [18, 19]. Therefore, the second model considers distinct coefficients for the characters at the beginning and the end of the string.

The nearest neighborhood of a character can also affect its retention time coefficient [19, 20]. The third model considers the influence of a character's direct neighborhood by considering coefficients for pairs of consecutive characters instead of coefficients for individual characters.

**Linear model:** Every character  $a \in \Sigma$  has a retention time coefficient  $t(a) \in \mathbb{Z}$ . The retention time of a string  $S$  is the sum of the retention time coefficients of its characters,

$$t_{\text{lin}}(S) := \sum_{i=1}^n t(a_i). \tag{1}$$

**Position-dependent model:** Characters at the first  $\gamma$  and the last  $\gamma$  positions of a string, where  $1 \leq \gamma \leq \lfloor \frac{n}{2} \rfloor$ , have distinct retention time coefficients. For  $i \leq \gamma$ , we denote the retention time coefficient of the  $i$ -th character by  $t_{\text{pre}}(a_i, i) \in \mathbb{Z}$  and the coefficient of the  $(n - i + 1)$ -th character by  $t_{\text{suf}}(a_{n-i+1}, i) \in \mathbb{Z}$ . The retention time of a string  $S$  is the sum of the corresponding retention time coefficients,

$$t_{\text{pos}}(S) := \sum_{i=1}^{\gamma} t_{\text{pre}}(a_i, i) + \sum_{j=\gamma+1}^{n-\gamma} t(a_j) + \sum_{k=1}^{\gamma} t_{\text{suf}}(a_{n-k+1}, k). \tag{2}$$

**Neighborhood-based model:** The model uses retention time coefficients  $t(a, b) \in \mathbb{Z}$  for pairs of characters  $a, b \in \Sigma$  that are consecutive in a given string  $S$ . The first and the last character  $a_1$  and  $a_n$  of  $S$  have additional coefficients  $t(-, a_1), t(a_n, -) \in \mathbb{Z}$ , as these characters have only one adjacent character in  $S$ . The retention time of  $S$  is the sum of all these coefficients,

$$t_{\text{nei}}(S) := t(-, a_1) + \left( \sum_{i=1}^{n-1} t(a_i, a_{i+1}) \right) + t(a_n, -). \tag{3}$$

The retention time coefficients for all three models can either be estimated from experimental data or taken from the literature. It is worth noting that the retention time coefficients might also be negative. Therefore, the retention time of a peptide does not depend linearly on the length of the peptide. We use a simple method for estimating the coefficients in the experimental evaluation and discuss limiting aspects of this method below.

### Problem definition

We recall the de novo peptide sequencing problem with respect to the symmetric difference scoring model [14]: Given a mass  $M$  and a set of fragment masses  $X$  (measured by the mass spectrometer), find a string  $S$  of mass  $M$  that minimizes  $|\text{TS}(S) \Delta X| = |\text{TS}(S) \setminus X| + |X \setminus \text{TS}(S)|$ . Equivalently to computing a string with mass  $M$  that minimizes  $|\text{TS}(S) \Delta X|$ , we can compute a string that maximizes  $|\text{TS}(S) \cap X| - |\text{TS}(S) \setminus X|$ , as  $X$  is a fixed input and  $S$  can be chosen. Throughout this paper, we assume that  $0, M \in X$ .

In this paper, we consider a variant of this problem that also considers the measured retention time  $T$  and a retention time prediction function  $t_* : \Sigma^* \rightarrow \mathbb{Z}$ . A function  $t_*(\cdot)$  can return negative values, as a substring can have a negative effect on the retention time of a string.

**Problem 1 (De Novo Sequencing Problem)** Let  $\Sigma$  be an alphabet of characters, with a mass  $m(a) \in \mathbb{N}$  for

each  $a \in \Sigma$ . Given a peptide mass  $M \in \mathbb{N}$ , a retention time  $T \in \mathbb{N}$ , a tolerance parameter  $\varepsilon \geq 0$  and a set  $X = \{x_i \in \mathbb{N} \mid i = 1, \dots, k\}$ , find a string  $S$  of characters in  $\Sigma$  with  $m(S) = M$  and  $|t(S) - T| \leq \varepsilon$  that minimizes  $|\text{TS}(S) \Delta X|$  among all strings with mass  $M$  and a retention time  $t_*(S) \in [T - \varepsilon, T + \varepsilon]$ .

**Methods**

**Algorithm for the symmetric difference scoring model**

We briefly describe the algorithm DeNovo $\Delta$  [14] for computing a string of mass  $M$  that minimizes  $|\text{TS}(S) \Delta X|$  without considering retention times. We refer to [14] for a detailed description and a proof of correctness. Then, we describe algorithms for solving the de novo sequencing problem for each considered prediction model.

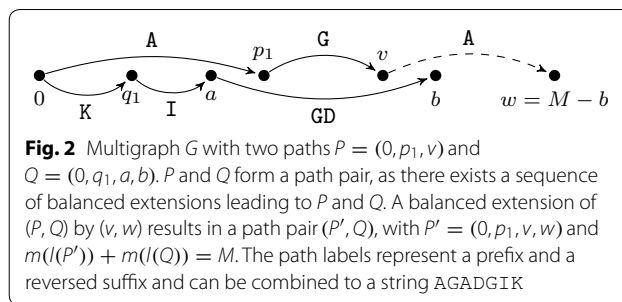
The search space of DeNovo $\Delta$  is modeled by a directed acyclic multigraph  $G = (V, E)$  based on the given set  $X$ . A vertex in  $G$  represents a mass and a path in  $G$  represents a string. For every mass  $m \in X$  there are two vertices  $m$  and  $M - m$  in  $G$ , i.e.  $V = \{m, M - m \mid m \in X\}$ . An edge in  $G$  is always directed from the smaller to the larger mass. Two vertices  $v$  and  $w$  are connected by an edge if there exists a string with mass  $w - v$ . For each such string with mass  $w - v$ , we add an edge from  $v$  to  $w$  to the multigraph and label it with this string. That is, if  $v$  and  $w$  are connected by an edge with label  $l(v, w)$ , there is also an edge from  $v$  to  $w$  for every permutation of  $l(v, w)$ . In practice, we only consider edges with a maximal label length  $p$ .

We denote the concatenation of the edge labels along a path  $P$  by  $l(P)$ . Let  $P = (0, v_1, \dots, v_k, M)$  be a path from vertex 0 to vertex  $M$ . Every traversed vertex  $v_i$  represents the mass of a prefix of the string  $l(P)$  and  $l(P)$  explains both  $v_i$  and  $M - v_i$  for every traversed vertex  $v_i$ .

The idea of DeNovo $\Delta$  for finding a string  $S$  of mass  $M$  that minimizes  $|\text{TS}(S) \Delta X|$  is to iteratively extend two paths both starting at vertex 0. One path represents a prefix and the other path a reversed suffix of  $S$ . DeNovo $\Delta$  extends both paths until the sum of their labels' masses is equal to  $M$  and then concatenates the prefix and the reversed suffix to a string of mass  $M$ .

**Definition 1 (Balanced extension)** Given two paths  $P$  and  $Q$  both starting at vertex 0, a *balanced extension* extends the path that represents the string of smaller mass by a single edge, unless the resulting paths represent strings with a total mass larger than  $M$ . An arbitrary path is extended if both paths represent strings with equal masses.

**Definition 2 (Path pair)** A *path pair* is a pair of paths  $P = (0, \dots, v)$  and  $Q = (0, \dots, a, b)$  in  $G$  that results from



a sequence of balanced extensions starting from two paths  $P_0 = (0)$  and  $Q_0 = (0)$ .

Figure 2 depicts an example of a path pair and a balanced extension. The set of masses that are explained by a path pair  $(P, Q)$  is the *partial theoretical spectrum*

$$\text{PTS}(P, Q, M) := \{ m(T), M - m(T) \mid T \in (\text{Pre}(l(P)) \cup \text{Pre}(l(Q))) \}. \tag{4}$$

The *score* of the path pair  $(P, Q)$  is the number of masses explained by the path pair that are in  $X$  minus the number of explained masses that are not in  $X$ , i.e.  $|\text{PTS}(P, Q, M) \cap X| - |\text{PTS}(P, Q, M) \setminus X|$ . The set of masses explained by an edge  $(v, w)$  is

$$\text{TSe}((v, w), M) := \{ m(T) + v, M - (m(T) + v) \mid T \in \text{Pre}(l(v, w)), m(T) \neq 0 \}. \tag{5}$$

**Lemma 1** For every path pair  $P = (0, \dots, v)$  and  $Q = (0, \dots, a, b)$  with  $v \leq b$  and  $v + b \leq M$  it holds that  $a \leq v \leq b$ . The balanced extension of  $(P, Q)$  by an edge  $(v, w)$  additionally explains all masses in  $N((v, w), (a, b)) = \text{TSe}((v, w), M) \setminus \text{TSe}((a, b), M)$ .

*Proof* Assume that there exists a path pair  $(P, Q)$  with  $v \leq a$ . This path pair results by definition from a sequence of balanced extensions. Consider the balanced extension in this sequence, where the last edge  $(a, b)$  of  $Q$  is added. In this step, either  $P$  ended in  $v$  or in some vertex  $v' < v$ . In both cases,  $a$  is the larger mass and  $Q$  represents the heavier string. Hence, the extension by  $(a, b)$  is not a balanced extension and  $(P, Q)$  is not a path pair.

Consider a balanced extension of  $(P, Q)$  by an edge  $(v, w)$ . The edge  $(v, w)$  explains all masses in  $\text{TSe}((v, w), M)$ . However, some of these masses might also be explained by  $(P, Q)$ . We show that  $\text{TSe}((v, w), M) \setminus \text{PTS}(P, Q, M) = N((v, w), (a, b))$ , i.e. that all masses explained by  $(v, w)$  that are also explained by  $(P, Q)$ , are explained by the last edge  $(a, b)$  of  $Q$ . We note that all masses in  $\text{TSe}((v, w), M)$  are larger than

$v$  and smaller than  $M - v$ . Moreover, all masses in PTS  $(P, Q, M)$  that are larger than  $v$  and smaller than  $M - v$  are explained by the edge  $(a, b)$ . Therefore, it follows that the balanced extension with  $(v, w)$  additionally explains all masses in  $N((v, w), (a, b))$ .  $\square$

Using Lemma 1, the algorithm DeNovo $\Delta$  [14] (Algorithm 1) computes a dynamic programming table  $DP$ . An entry  $DP[v, (a, b)]$  contains the optimal score of a path pair ending at the vertex  $v$ , respectively at the edge  $(a, b)$ . As a base case, we add a loop edge  $(0, 0)$  to the graph and initialize  $DP[0, (0, 0)] = 2$ , because the path pair representing two empty strings explains the masses  $0, M \in X$ . Given the optimal score  $DP[v, (a, b)]$ , the algorithm considers all possible balanced extensions of the corresponding path pair with outgoing edges of  $v$ . By Lemma 1, the additionally explained masses of such a balanced extension can be computed only given the last vertex  $v$  and the last edge  $(a, b)$  of the two paths. The score of the resulting new path pair can be computed by adding

$$\text{gain}((v, w), (a, b)) := |N((v, w), (a, b)) \cap X| - |N((v, w), (a, b)) \setminus X| \quad (6)$$

```

1   $DP[v, (a, b)] \leftarrow -\infty$  for all  $(a, b) \in E$  and all  $v \in V$ 
2   $DP[0, (0, 0)] \leftarrow 2$ 
3  for  $v \in V$  in ascending order do
4      foreach  $(a, b) \in E$  in lexicograph. asc. order with  $DP[v, (a, b)] \neq -\infty$  do
5          foreach  $(v, w) \in E$  with  $w + b \leq M$  do
6              if  $w \leq b$  then
7                   $DP[w, (a, b)] \leftarrow$ 
8                       $\max ( DP[w, (a, b)], DP[v, (a, b)] + \text{gain}((v, w), (a, b)) )$ 
9              else
10                  $DP[b, (v, w)] \leftarrow$ 
11                      $\max ( DP[b, (v, w)], DP[v, (a, b)] + \text{gain}((v, w), (a, b)) )$ 
12             end
13         end
14     end

```

to the score  $DP[v, (a, b)]$ . The table entry of the new path pair is updated if the new score exceeds the value stored in this entry at this step of the algorithm. The optimal score for a string of mass  $M$  is equal to the maximum value of an entry  $DP[M - b, (a, b)]$  among all edges  $(a, b)$  in  $G$ . A path pair with this score can be reconstructed starting from this entry. The combination of the corresponding prefix and reversed suffix then leads to the desired string of mass  $M$ . The time complexity of DeNovo $\Delta$  is  $\mathcal{O}(|V| \cdot |E| \cdot d \cdot p)$ , where  $d$  is the maximal out-degree of a vertex in  $G$  and  $p$  is the maximal length of an edge label [14].

#### Algorithm for the linear prediction model

In the following subsections, we develop an algorithm for the de novo sequencing problem (Problem 1). We have to consider three aspects when taking into account the retention time information. First, we need to define the predicted retention time of a path pair in  $G$ . Second, we have to compute the effect of a balanced extension on the predicted retention time of a path pair. Third, we need to find optimal substructures of paths from 0 to  $M$  in  $G$  with an optimal score and a feasible predicted retention time.

Algorithm 1: DeNovo $\Delta$  [14].

In this subsection, we consider the linear retention time prediction model. We note that the retention time of a path pair  $P = (0, \dots, v)$  and  $Q = (0, \dots, a, b)$  with  $a \leq v \leq b$  is the sum of the retention times of both substrings  $t = t_{\text{lin}}(\perp(P)) + t_{\text{lin}}(\perp(Q))$ . Moreover, the retention time  $t'$  of a path pair obtained from  $(P, Q)$  by applying a balanced extension by some edge  $(v, w)$  can be computed as  $t' = t + t_{\text{lin}}(\perp(v, w))$ . That is, we only need  $t$  and the edge label  $l(v, w)$  for computing  $t'$ .

However, it is not sufficient to only store the optimal score  $DP[v, (a, b)]$  of any path pair ending in  $v$ , respectively  $(a, b)$ , and its retention time to compute a solution for our problem. There can be multiple path pairs ending in the same vertex and the same edge with different retention times. If we consider an optimal solution and its sequence of path pairs computed by the algorithm, a path pair  $P = (0, \dots, v)$  and  $Q = (0, \dots, a, b)$  in this sequence does not necessarily have an optimal score among all path pairs ending in  $v$  and  $(a, b)$ . Nevertheless, its score is optimal among all path pairs with the same retention time that end in  $v$  and  $(a, b)$ . Therefore, we need to store for each possible retention time  $t$  the optimal score of a path pair ending in vertex  $v$  and edge  $(a, b)$ .

DeNovo $\Delta$ Lin (Algorithm 2) stores for each entry  $DP[v, (a, b)]$  an array containing a score for every possible retention time  $t$ .  $DP[v, (a, b)][t]$  is the optimal score for a path pair ending in  $v$ , respectively  $(a, b)$ , with retention time  $t$ . For a given vertex  $v$  and an edge  $(a, b)$ , the algorithm performs balanced extensions by all outgoing edges  $(v, w)$  of  $v$ . For every balanced extension and every feasible retention time  $t$ , the algorithm then computes the new retention time  $t'$  and the new score of the resulting path pair and updates the corresponding entry in the table. We can see by an inductive argument that the optimal scores in the table are computed correctly. As the base case, we note that  $DP[0, (0, 0)][0] = 2$  is correct, as an empty path pair explains the masses  $\{0, M\} \subseteq X$  and has retention time 0. As soon as the entry  $DP[v, (a, b)]$  is reached in line 7, all optimal scores for path pairs ending in vertex  $v$  and edge  $(a, b)$  have been computed. This holds by induction, as every possible balanced extension leading to a path pair ending in  $v$  and  $(a, b)$  has already been considered (given the optimal score of a preceding path pair). Moreover, the array in  $DP[v, (a, b)]$  is not further modified as soon as the algorithm reaches the vertex  $v$  and the edge  $(a, b)$  in line 7. Therefore, the invariant holds that, if the algorithm considers a vertex  $v$  and an edge  $(a, b)$  in line 7, the corresponding entry  $DP[v, (a, b)]$  contains the optimal score for each feasible retention time.

After computing all entries  $DP[v, (a, b)]$ , the optimal score of a string with retention time  $t$  is

$\max_{(a,b) \in E} DP[M - b, (a, b)][t]$ . We are interested in optimal strings with a predicted retention time  $t \pm \varepsilon$ . Therefore, we iterate over all entries  $DP[M - b, (a, b)][t]$  for  $(a, b) \in E$  and all feasible retention times  $t \in [T - \varepsilon, T + \varepsilon]$  to find the optimal score of a string with a feasible predicted retention time. We can reconstruct a corresponding string starting from the corresponding entry in  $DP$ .

The running time of DeNovo $\Delta$  is in  $\mathcal{O}(|V| \cdot |E| \cdot d \cdot p)$  [14], where  $d$  is the maximal out-degree of a vertex in  $G$  and  $p$  is the maximal length of an edge label. The additional overhead of DeNovo $\Delta$ Lin (loop starting at line 8 in Algorithm 2) is to iterate over all feasible retention times  $t$  for each entry  $DP[v, (a, b)]$  and compute the new retention time  $t'$ .

The number of scores to be stored varies depending on the entry and the retention time coefficients. For a path pair ending in  $v$ , respectively  $(a, b)$ , we have to consider all retention times in  $[rt_{\text{min}} \cdot (v + b), rt_{\text{max}} \cdot (v + b)]$ , where  $rt_{\text{min}}$  and  $rt_{\text{max}}$  are the minimum and the maximum retention time per mass unit. For example, we only store one optimal score in entry  $DP[0, (0, 0)]$ , but up to  $\lceil rt_{\text{max}} \cdot M - rt_{\text{min}} \cdot M \rceil$  scores in entries  $DP[M - b, (a, b)]$  for  $(a, b) \in E$ . The time complexity of DeNovo $\Delta$ Lin is in  $\mathcal{O}(|V| \cdot |E| \cdot |RT_M| \cdot d \cdot p)$ , where  $|RT_M|$  denotes the number of possible retention times for a string of mass  $M$ . In practice, most entries  $DP[v, (a, b)]$  contain only few scores, as we only store the score for a retention time  $t$  if there is a path pair ending in  $v$  and  $(a, b)$  with predicted retention time  $t$ . Therefore, it is advisable to use a memory-efficient data structure instead of an array to reduce the memory consumption of the algorithm.

This approach is flexible and can be extended to compute suboptimal solutions, e.g. the  $k$  best-scoring strings, using similar techniques as described in [14]. The implementation of this algorithm supports computing both the best and the  $k$  best strings for a given input.

#### Algorithm for the position-dependent prediction model

In the position-dependent prediction model, the retention time of a string  $S$  is not equal to the retention time of all permutations of  $S$ . This is due to the fact that the retention time coefficient of a character in the first and the last  $\gamma$  positions of the string may be different from the coefficient of the same character at another position. Therefore, we have to distinguish the prefix and the suffix path of a path pair  $(P, Q)$ , with  $P = (0, \dots, v)$ ,  $Q = (0, \dots, a, b)$ , and  $a \leq v \leq b$ , in order to compute its predicted retention time. This was not necessary for DeNovo $\Delta$  and DeNovo $\Delta$ Lin, as both the score and the

predicted retention time (in the linear prediction model) do not depend on which of the two paths represents the prefix.

Let us assume that  $P$  is the prefix path and  $Q$  is the suffix path of a path pair  $(P, Q)$ . We compute the retention time of  $(P, Q)$  by summing the retention times  $t_P$  and  $t_Q$  of the path labels,

$$t_P := \sum_{a_i \in l(P)} \begin{cases} t_{\text{pre}}(a_i, i) & i \leq \gamma \\ t(a_i) & i > \gamma \end{cases} \quad (7)$$

$$t_Q := \sum_{a_j \in l(Q)} \begin{cases} t_{\text{suf}}(a_j, j) & j \leq \gamma \\ t(a_j) & j > \gamma. \end{cases}$$

If we want to update the retention time after a balanced extension of  $(P, Q)$  by an edge  $(v, w)$ , we have to compute the retention time of the edge label  $l(v, w)$ . This retention time depends on whether the edge label contains some of the first or the last  $\gamma$  characters of a solution string  $S$  of mass  $M$ . However, there can be multiple such solution strings resulting from different further balanced extensions of this path pair.

```

1  foreach  $(a, b) \in E$  and  $v \in V$  do
2  |    $DP[v, (a, b)] \leftarrow$  array with entries  $-\infty$  for each feasible retention time  $t$ 
3  end
4   $DP[0, (0, 0)][0] \leftarrow 2$ 
5  for  $v \in V$  in ascending order do
6  |   foreach  $(a, b) \in E$  in lexicograph. asc. order with  $a \leq v \leq b$  do
7  |   |   foreach  $(v, w) \in E$  with  $w + b \leq M$  do
8  |   |   |   foreach entry  $t$  in  $DP[v, (a, b)]$  do
9  |   |   |   |    $t' \leftarrow t + t_{\text{lin}}(l(v, w))$ 
10  |   |   |   |   if  $w \leq b$  then
11  |   |   |   |   |    $DP[w, (a, b)][t'] \leftarrow$ 
12  |   |   |   |   |    $\max ( DP[w, (a, b)][t'], DP[v, (a, b)][t] + \text{gain}((v, w), (a, b)) )$ 
13  |   |   |   |   |   else
14  |   |   |   |   |   |    $DP[b, (v, w)][t'] \leftarrow$ 
15  |   |   |   |   |   |    $\max ( DP[b, (v, w)][t'], DP[v, (a, b)][t] + \text{gain}((v, w), (a, b)) )$ 
16  |   |   |   |   |   end
17  |   |   |   end
18  |   |   end
19  |   end
20  end

```

Algorithm 2: DeNovo $\Delta$ Lin– Linear retention time prediction model.

We can decide whether  $l(v, w)$  contains some of the first  $\gamma$  characters given the length  $k$  of  $l(P)$  without knowing the solution string  $S$ . If  $k \geq \gamma$ , the edge label clearly does not contain any of the first  $\gamma$  characters of any solution resulting from extending  $(P, Q)$ . Likewise, we know that  $l(v, w)$  contains none of the  $\gamma$  last characters if  $l(Q)$  has more than  $\gamma$  characters. However, if  $l(Q)$  has less than  $\gamma$  characters, we cannot decide whether  $l(v, w)$  contains some of the last  $\gamma$  characters without knowing the length of  $S$ .

Let us assume for now that  $l(v, w)$  does not contain some of the last  $\gamma$  characters of the solution. The retention time of the new path pair resulting from the balanced extension of  $(P, Q)$  by the edge  $(v, w)$  is

$$t' = t + \sum_{a_i \in l(v, w)} \begin{cases} t_{\text{pre}}(a_i, i) & i + k \leq \gamma \\ t(a_i) & i + k > \gamma. \end{cases} \quad (8)$$

If  $P$  would be the suffix path,  $t_{\text{pre}}(a_i, i)$  would be replaced by  $t_{\text{suf}}(a_i, i)$  in the above equation.

It is important that the above assumption holds for every balanced extension leading to a solution string  $S$ . Otherwise, the retention time of the new path pair is not computed correctly. We cannot check if our assumption holds while computing the new retention time after a balanced extension. However, given a solution string  $S$  and a path pair that represents a prefix and a suffix of  $S$ , we can check if either the balanced extension leading to this path pair or a preceding balanced extension does not satisfy the assumption. If so, either the prefix or the suffix path label has at least  $n - \gamma$  characters, where  $n$  is the length of  $S$ . This also holds for all subsequent path pairs, as we only add characters to path labels in a balanced extension.

Therefore, when reconstructing a solution from the dynamic programming table, we have to additionally check, if one of the path labels has  $n - \gamma$  or more characters, before we combine them to a solution string. If so, the assumption was not fulfilled at some step and we discard this solution, as its retention time was not computed correctly. Note that we cannot consider these strings, unless they can be constructed by another sequence of balanced extensions. However, it is very unlikely that the assumption is not fulfilled in practice, as we consider small values of  $\gamma$ . We never observed such a situation in our evaluation using  $\gamma = 2$ .

Given the sequence of path pairs of an optimal solution, a path pair in this sequence has an optimal score among all path pairs with the same retention time. However, we have to store some additional information to compute a solution with respect to the position-dependent prediction model. First, we have to store whether  $P$  is a prefix or a suffix path. Second, we have to store the length of both path labels, unless they are larger than  $\gamma$ .

DeNovo $\Delta$ Pos (Algorithm 3) stores the optimal scores of path pairs ending in  $v$  and  $(a, b)$  in an array with an entry for every retention time  $t$ , the lengths  $\alpha$  and  $\beta$  of the path labels and a Boolean variable  $pre$  indicating if the path ending in  $v$  is the prefix or the suffix path. We store the length of the path labels only up to length  $\gamma$ , as the exact length is only important as long as the path labels have less than  $\gamma$  characters.

If the algorithm reaches an entry  $DP[v, (a, b)]$  in line 7, all optimal scores for path pairs ending in vertex  $v$  and edge  $(a, b)$  have been computed correctly, as all balanced extensions leading to such path pairs have already been considered. Given the optimal score of a path pair, the algorithm performs every possible balanced extension with an outgoing edge of  $v$ , computes the new score and retention time, and updates the corresponding entries.

We reconstruct a solution starting from a path pair ending in some vertex  $M - b$  and some edge  $(a, b)$ . The algorithm additionally verifies that both the prefix and the suffix path label have more than  $\gamma$  characters. DeNovo $\Delta$ Pos considers at most  $2 \cdot \gamma^2 \cdot |RT_M|$  optimal scores for each table entry  $DP[v, (a, b)]$ , where  $|RT_M|$  is the number of possible retention times for a string of mass  $M$ . Therefore, the running time is in  $\mathcal{O}(|V| \cdot |E| \cdot |RT_M| \cdot \gamma^2 \cdot d \cdot p)$ , where  $d$  is the maximal out-degree of a vertex in  $G$  and  $p$  is the maximal length of an edge label.

#### Algorithm for the neighborhood-based prediction model

The neighborhood-based model predicts the retention time of a string  $S$  by considering all pairs of consecutive characters. We define the predicted retention time of a path pair  $(P, Q)$  as follows. The retention time of the path label  $l(P)$  is the sum of the retention time coefficients of the pairs of consecutive characters and the additional coefficient of the first character. Note that we consider only one coefficient for the last character in the prefix, as the other coefficient depends on the next balanced extension or the last character of  $l(Q)$ . The retention time of  $l(Q)$  is defined analogously considering that the  $l(Q)$  is a reversed suffix of the solution string  $S$ . We compute the retention time of  $(P, Q)$  by summing the retention times of both path labels (Fig. 3). That is, the retention time of  $(P, Q)$  is

$$t_{\text{nei}}(P, Q) := t(-, p_1) + \left( \sum_{i=1}^{n-1} t(p_i, p_{i+1}) \right) + \left( \sum_{i=m}^2 t(q_i, q_{i-1}) \right) + t(q_1, -), \quad (9)$$

where  $l(P) = p_1, \dots, p_n$  and  $l(Q) = q_1, \dots, q_m$  are the path labels of  $(P, Q)$ .



```

1  foreach  $(a, b) \in E$  and  $v \in V$  do
2  |    $DP[v, (a, b)] \leftarrow (|RT_M| \times \gamma \times \gamma \times 2)$ -array with entries  $-\infty$ 
3  end
4   $DP[0, (0, 0)][0, 0, 0, 0] \leftarrow 2$ 
5  for  $v \in V$  in ascending order do
6  |   foreach  $(a, b) \in E$  in lexicograph. asc. order with  $a \leq v \leq b$  do
7  |   |   foreach  $(v, w) \in E$  with  $w + b \leq M$  do
8  |   |   |   foreach entry  $(t, \alpha, \beta, pre)$  in  $DP[v, (a, b)]$  do
9  |   |   |   |    $t' \leftarrow$  retention time of resulting path pair
10 |   |   |   |   if  $pre$  then
11 |   |   |   |   |    $\alpha' \leftarrow \max(\gamma, \alpha + |l(v, w)|)$ ;  $\beta' \leftarrow \beta$ 
12 |   |   |   |   else
13 |   |   |   |   |    $\alpha' \leftarrow \alpha$ ;  $\beta' \leftarrow \max(\gamma, \beta + |l(v, w)|)$ 
14 |   |   |   |   end
15 |   |   |   |   if  $w \leq b$  then
16 |   |   |   |   |    $DP[w, (a, b)][t', \alpha', \beta', pre] \leftarrow$ 
17 |   |   |   |   |    $\max ( DP[w, (a, b)][t', \alpha', \beta', pre], DP[v, (a, b)][t, \alpha, \beta, pre]+$ 
18 |   |   |   |   |    $\text{gain}((v, w), (a, b)) )$ 
19 |   |   |   |   else
20 |   |   |   |   |    $DP[b, (v, w)][t', \alpha', \beta', \neg pre] \leftarrow$ 
21 |   |   |   |   |    $\max ( DP[b, (v, w)][t', \alpha', \beta', \neg pre], DP[v, (a, b)][t, \alpha, \beta, pre]+$ 
22 |   |   |   |   |    $\text{gain}((v, w), (a, b)) )$ 
23 |   |   |   |   end
24 |   |   |   end
25 |   |   end
26 |   end
27 end

```

Algorithm 3: DeNovo $\Delta$ Pos – Position-dependent prediction model.



**Fig. 3** The retention time  $t$  of a path pair  $(P, Q)$  is the sum of the retention time coefficients up to the last characters  $p_2$  and  $q_3$ . The path pair  $(P', Q)$  resulting from a balanced extension of  $(P, Q)$  by an edge with label  $l_1 l_2$  has retention time  $t + t(p_2, l_1) + t(l_1, l_2)$ . A path pair  $(P', Q)$  with  $m(l(P')) + m(l(Q)) = M$  can be combined to a solution string  $S$  by concatenating  $l(P')$  and the reversed string of  $l(Q)$ . The retention time of  $S$  is  $t_{nei}(P', Q) + t(l_2, q_3)$

We can update the retention time after a balanced extensions of  $(P, Q)$  as follows. Consider a balanced extension of the prefix path  $P$  by an edge  $(v, w)$  with  $l(v, w) = l_1 \dots l_k$ . Let  $p_n$  be the last character of  $l(P)$ . The retention time  $t'$  of the new path pair resulting from the balanced extension is

$$t' = t_{nei}(P, Q) + t(p_n, l_1) + \sum_{i=1}^{k-1} t(l_i, l_{i+1}). \quad (10)$$

The retention time after a balanced extension of the suffix path  $Q$  is defined analogously (again considering the  $l(Q)$  is a reversed suffix).

Note that the retention time of a solution  $S$  is not the sum of the retention times of a prefix of  $S$  and its complementary suffix. We have to additionally consider the coefficient of the last character of the prefix and the first

character of the suffix, which are consecutive in  $S$ . If we combine the path labels of a path pair  $(P', Q)$  to a string  $S$  (Fig. 3), the retention time of  $S$  is  $t_{nei}(P', Q) + t(p_n, q_m)$ , where  $p_n$  and  $q_m$  are the last characters of the prefix  $l(P)$  and the reversed suffix  $l(Q)$ .

DeNovoΔNei (Algorithm 4) stores for every path pair  $(P, Q)$  ending in vertex  $v$  and edge  $(a, b)$  the optimal score for each retention time  $t$ , last character  $p$  of the path ending in  $v$ , and a Boolean variable *pre* indicating if  $P$  is the prefix path. As a base case, the algorithm stores the optimal score for a path pair ending in vertex 0 and the loop edge  $(0, 0)$  as  $DP[0, (0, 0)][0, -, 0] = 2$ . The algorithm considers the vertices and edges of  $G$  in ascending order. After considering all possible path pairs, the optimal score can be computed by considering all entries  $DP[M - b, (a, b)]$  and the feasible solutions for path pairs ending in these vertices and edges.

```

1  foreach  $(a, b) \in E$  and  $v \in V$  do
2  |    $DP[v, (a, b)] \leftarrow (|RT_M| \times |\Sigma| \times 2)$ -array with entries  $-\infty$ 
3  end
4   $DP[0, (0, 0)][0, -, 0] \leftarrow 2$ 
5  for  $v \in V$  in ascending order do
6  |   foreach  $(a, b) \in E$  in lexicograph. asc. order with  $a \leq v \leq b$  do
7  |   |   foreach  $(v, w) \in E$  with  $w + b \leq M$  do
8  |   |   |   foreach entry  $(t, \mathbf{p}, pre)$  in  $DP[v, (a, b)]$  do
9  |   |   |   |    $t' \leftarrow$  retention time of resulting path pair
10 |   |   |   |   if  $w \leq b$  then
11 |   |   |   |   |    $\mathbf{p}' \leftarrow$  last character of  $l(v, w)$ 
12 |   |   |   |   |    $DP[w, (a, b)][t', \mathbf{p}', pre] \leftarrow$ 
13 |   |   |   |   |    $\max ( DP[w, (a, b)][t', \mathbf{p}', pre], DP[v, (a, b)][t, \mathbf{p}, pre] +$ 
14 |   |   |   |   |    $\text{gain}((v, w), (a, b)) )$ 
15 |   |   |   |   else
16 |   |   |   |   |    $\mathbf{p}' \leftarrow$  last character of  $l(a, b)$ 
17 |   |   |   |   |    $DP[b, (v, w)][t', \mathbf{p}', -pre] \leftarrow$ 
18 |   |   |   |   |    $\max ( DP[b, (v, w)][t', \mathbf{p}', -pre], DP[v, (a, b)][t, \mathbf{p}, pre] +$ 
19 |   |   |   |   |    $\text{gain}((v, w), (a, b)) )$ 
20 |   |   |   |   end
21 |   |   |   end
22 |   |   end
23 |   end
24 end

```

Algorithm 4: DeNovo $\Delta$ Nei – Neighborhood-based prediction model.

The algorithm considers at most  $2 \cdot |\Sigma| \cdot |RT_M|$  optimal scores for each pair of a vertex  $v$  and an edge  $(a, b)$ , where  $|RT_M|$  is the number of possible retention times for a string of mass  $M$  and  $|\Sigma|$  is the size of the considered alphabet. The running time of DeNovo $\Delta$ Nei is in  $\mathcal{O}(|V| \cdot |E| \cdot |RT_M| \cdot |\Sigma| \cdot d \cdot p)$ , where  $d$  is the maximal out-degree of a vertex,  $p$  is the maximal length of an edge label, and  $|RT_M|$  is the number of feasible retention times for a string of mass  $M$ .

### Experimental evaluation

In this section, we study the performance of our algorithms for de novo peptide sequencing with retention time prediction. In our evaluation, we want to clearly expose the effect of considering the retention time information rather than studying the identification rates compared to state-of-the-art de novo sequencing software, such as UniNovo [6] or Novor [8]. We compare the identification rates of the proposed algorithms with the identification rates of DeNovo $\Delta$  [14], as this algorithm uses the same symmetric difference scoring model, while other available tools use different scoring models. Note that we use a very simple scoring function that only considers if a mass has been measured by the instrument, but no other information, such as the intensity of the signal. While this is sufficient for studying the effect of considering the retention time information, such a scoring function is typically not suitable for real applications. However, our algorithms can support more sophisticated scoring models that also take into account the signal intensities measured by the mass spectrometer. We refer to [14] for one example of such a scoring function that is supported by the current implementation of our algorithms.

We first describe the considered dataset and a method for estimating the parameters of the three models. Then, we compare the identification rates of the proposed algorithms to the identification rate of DeNovo $\Delta$  [14].

### Dataset

We use the SWATH-MS Gold Standard (SGS) dataset (<http://www.peptideatlas.org>, identifier PASS00289, [15]) with measurements of 422 synthesized peptides. Specifically, we consider the 944 spectra of synthesized peptides from DDA-experiments that have also been considered in [14]. The raw profile spectra were centroided (peak-picked) using the tool qtofpeak-picker [21]. The spectra have been analyzed using the database search tool Comet [22] using the very restricted database containing only the 422 synthesized peptides. In our evaluation, we only considered spectra from doubly-charged peptides (as reported by Comet) and assumed that all measured

fragment masses are singly charged. Peptideprophet [23] has been used to validate the results.

We used the sequences identified by Comet as gold standard and considered a peptide to be identified by one of the considered algorithm, if the exact sequence has been computed as the best-scoring solution, respectively one of the 5, 10, or 100 best-scoring solutions.

### Retention time coefficient estimation

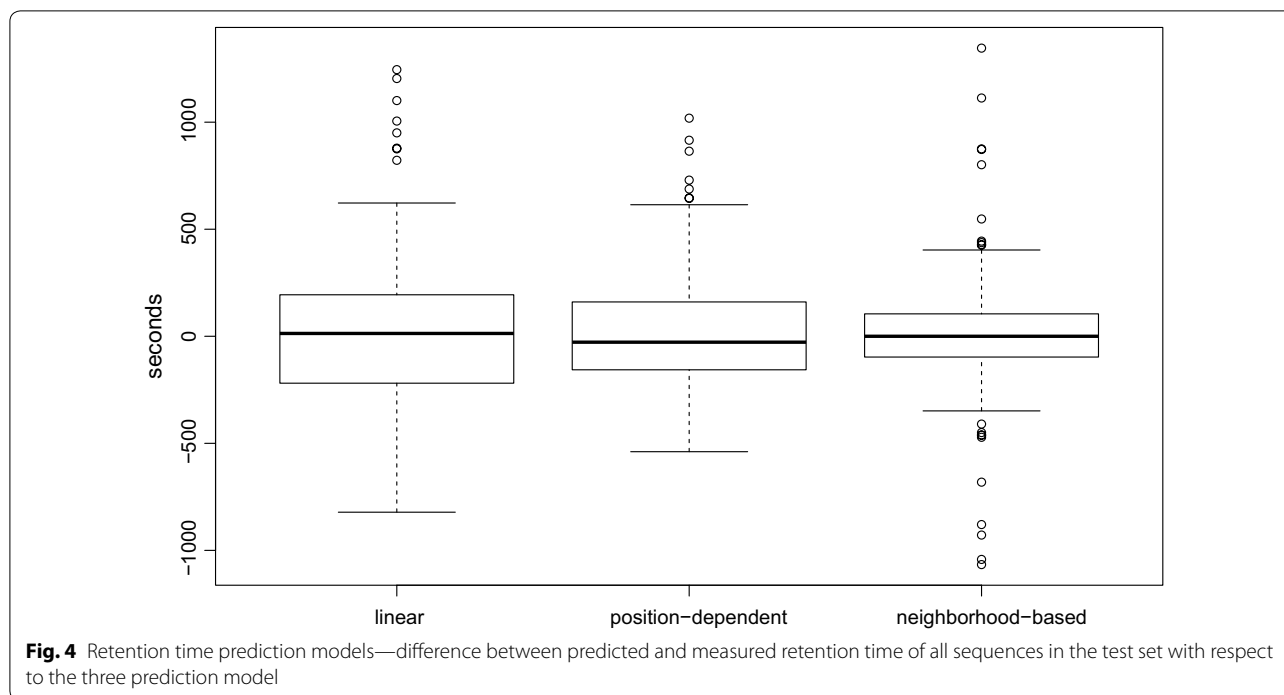
In this work, we are mainly interested in the algorithmic problem of using retention time information for de novo sequencing and do not focus on efficient procedures for estimating the coefficients of retention time prediction models. We use linear regression for estimating the coefficients for our three retention time models.

We randomly split the 944 spectra into a training set with 80% of the spectra (755 spectra) and a test set with the remaining 20% of the spectra (189 spectra). We use the training set to estimate the retention time coefficients and the test set to select a tolerance parameter  $\varepsilon$ . In a linear regression approach, we choose the coefficients such that the sum of the squared loss  $\sum_{S,T} (T - t(S))^2$  is minimized, where  $T$  is the measured retention time, and  $t(S)$  the predicted retention time of the sequence  $S$ .

For example, we estimate the coefficients of the linear model by first computing the character frequency vector for each string in the dataset. The character frequency vector of a string is a vector of length  $|\Sigma|$  that indicates how often a character occurs in the string. For example, the occurrence vector of the string AGA has value 2 at entry A, value 1 at entry G and value 0 at all other entries. Then, the retention time of a string  $S$  is the scalar product of its character frequency vector  $freq(S)$  and the vector of the retention time coefficients  $ct$ . Standard software tools for statistical methods [24] can be used to compute  $ct$ , such that  $\sum_i (T_i - \langle ct, freq(S) \rangle)^2$  is minimized.

We chose the tolerance parameter  $\varepsilon$  independently for each prediction model by considering the difference between the measured and the predicted retention time of the sequences in the test set. Figure 4 shows the differences between the predicted and the measured retention times for all three models on the test dataset. We set  $\varepsilon$  to half the difference between the maximum error  $e_{\max}$  and the minimum error  $e_{\min}$ , i.e.  $\varepsilon = (e_{\max} - e_{\min})/2$ . Concretely, we set  $\varepsilon = 1000$  seconds for the linear prediction model and  $\varepsilon = 750$  seconds for the position-dependent model.

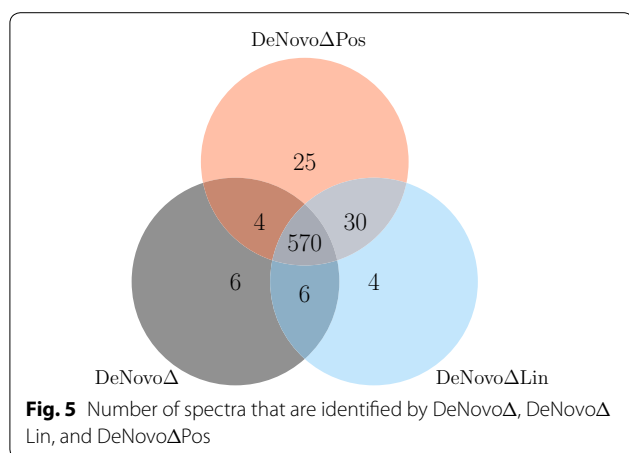
The neighborhood-based prediction model has a very large predictive error for several sequences due to the small training dataset. Several coefficients are estimated based on few observations and others cannot be estimated at all. Therefore, we cannot extensively evaluate



the identification rates of our algorithm with the neighborhood-based prediction model, as a much larger training dataset for estimating all parameters would be necessary. For our limited evaluation, we ignore the 5 largest and the 5 smallest retention time errors when picking the tolerance parameter and use  $\epsilon = 500$  seconds.

**Comparison of DeNovo $\Delta$ Lin and DeNovo $\Delta$ Pos**

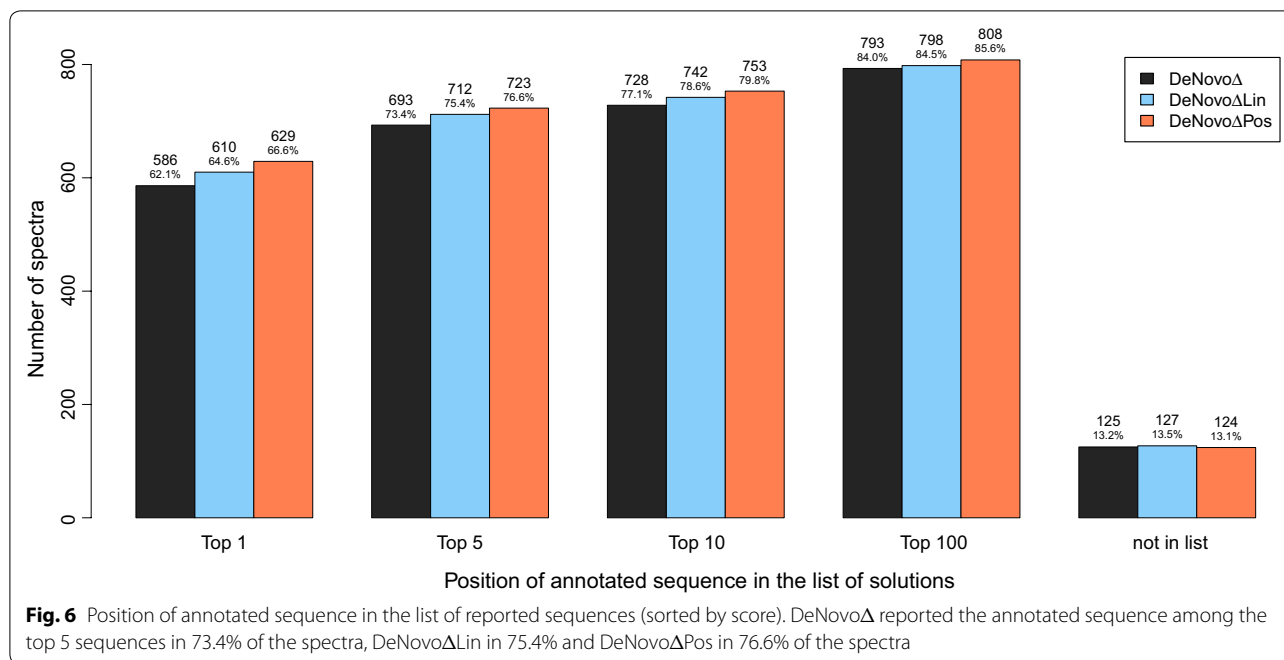
We analyzed the 944 considered spectra with DeNovo $\Delta$ Lin and DeNovo $\Delta$ Pos. Both algorithms compute all solutions with a score of at least 90% of the optimal score and a predicted retention time within the tolerance range. Figure 5 shows the number of annotated sequences



reported as best-scoring sequence by the three considered algorithms. While the majority of the spectra are either identified by all algorithms or not at all, 59 spectra are only identified when considering the retention time information.

Figure 6 shows a comparison of the identification rates with respect to the 5, 10, and 100 best-scoring sequences of DeNovo $\Delta$  [14], DeNovo $\Delta$ Lin, and DeNovo $\Delta$ Pos. Without considering the retention time, DeNovo $\Delta$  reported the annotated sequence as best-scoring sequence for 586 spectra (62.1%). Considering the linear retention time prediction model, DeNovo $\Delta$ Lin computed the annotated sequence with an optimal score for 610 spectra (64.6%). DeNovo $\Delta$ Pos considers the position-dependent prediction model and achieved the highest identification rate. The annotated sequence was reported as best-scoring sequence for 629 spectra (66.6%). The performance improvement decreases with increasing number of considered candidate sequences.

However, a filtering approach that considers the top 100 sequences reported by DeNovo $\Delta$ , would not be as successful as the proposed algorithms. While the annotated sequence was reported by DeNovo $\Delta$  for 793 spectra among the top 100 sequences, DeNovo $\Delta$ Lin reported it in 798 cases and DeNovo $\Delta$ Pos in 808 cases. Even an optimal filtering approach by retention time would miss the sequences that have not been reported by DeNovo $\Delta$ . For six spectra, DeNovo $\Delta$ Lin and DeNovo $\Delta$ Pos did not report the annotated sequence, where DeNovo $\Delta$  did



report it, as the predicted retention time of the annotated sequence was not in the chosen tolerance range.

The length of a peptide affects its retention time. However, the considered prediction models do not take into account the peptide's length and use the same coefficients for all peptide lengths. There is not necessarily a linear correlation between the length of a peptide and its retention time, as the coefficients can be positive or negative. Our models do not perform equally well on short and long peptides. Figure 7 shows a distribution of the number of identified spectra with respect to the length of the corresponding peptide sequence. DeNovoΔPos shows the best performance for peptides with fewer than 14 amino acids. For longer peptides, the linear prediction model shows a superior identification rate on the considered dataset.

## Discussion and conclusion

### Discussion

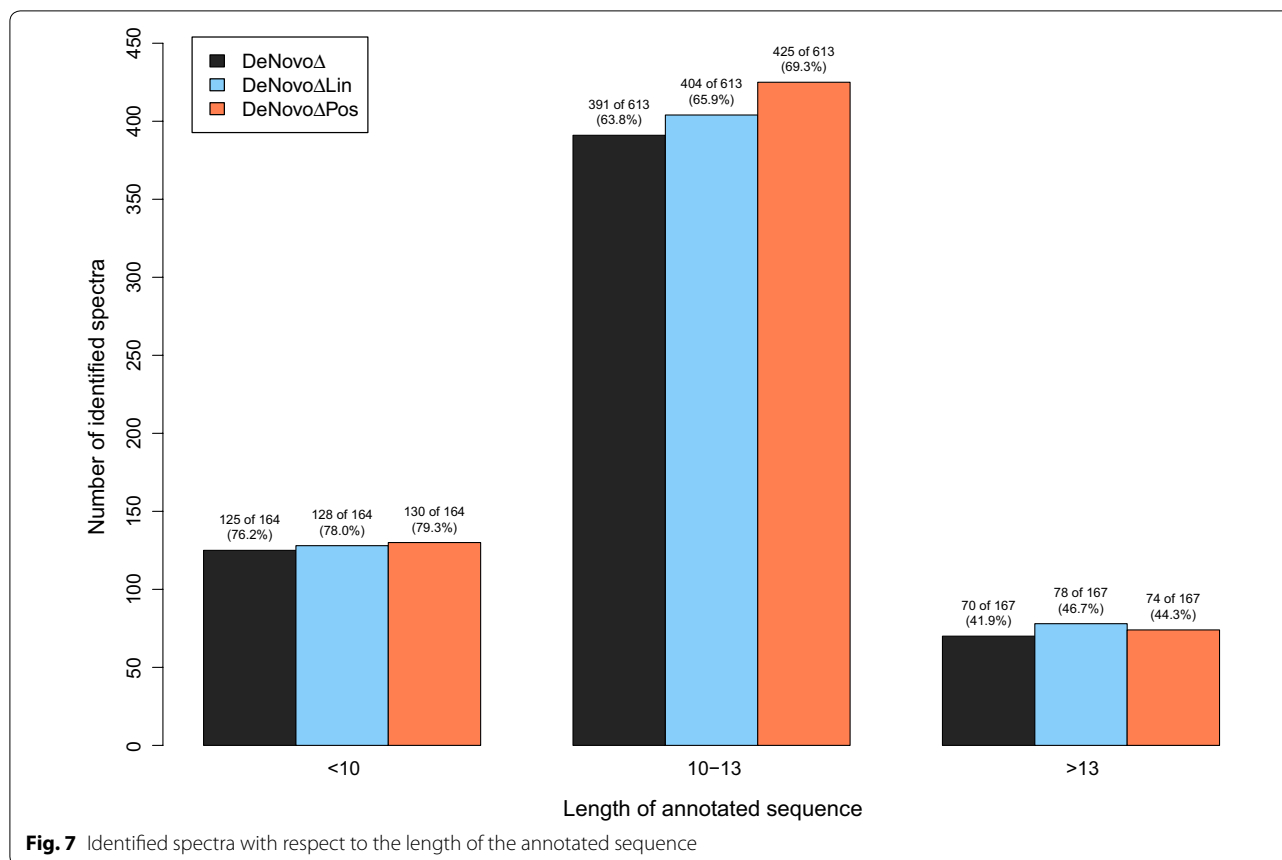
An accurate retention time prediction model is crucial for exploiting the retention time information successfully. The identification rates of our algorithms depend on the choice of the tolerance parameter  $\varepsilon$ . Increasing  $\varepsilon$  diminishes the effect of considering the retention time, while decreasing  $\varepsilon$  might exclude the correct sequence from the search space.

In our evaluation, we considered a limited training dataset for estimating the retention time coefficients. While we have to estimate a small set of coefficients for our linear prediction model, the neighborhood-based

prediction model has many retention time coefficients. Estimating these coefficients requires a large training dataset, as each coefficient needs to be estimated based on a sufficiently large set of observations. A much larger training set would be necessary to get a robust estimate of the retention time coefficients for this model. Our models fail to predict the retention time of some sequences accurately considering the available training data. To avoid excluding the correct sequence from the search space, we had to choose large tolerance parameters. By improving the predictive power of the models, e.g. using a larger training set or a more sophisticated parameter estimation, the tolerance parameter can be decreased, which increases the identification rates of our algorithms.

To get a glimpse on the performance of DeNovoΔNei, we set  $\varepsilon = 500$  (in seconds) and analyzed the spectra from the test set, where the correct sequence was not excluded due to the predictive error. In three cases, the annotated sequence was reported by DeNovoΔNei, but by no other considered algorithm. The position of the annotated sequence improved compared to the position reported by DeNovoΔPos for 12 spectra.

Our prediction models do not consider several other properties of a peptide that affect its retention time. For example, the length of a peptide has an influence on its retention time. More evolved prediction models [18, 19] integrate a correction for the peptide length. The prediction models considered in this work cannot account for the peptide length. However, as suggested in [19], a separate set of retention time coefficient can be estimated for



short peptides in order to improve the prediction accuracy. This approach needs an even larger training dataset in order to accurately estimate the coefficients.

The running time of our prototypical implementations is in some cases not yet practical. DeNovoΔLin needs less than 3 seconds per spectra for half of the considered spectra, but several hours in exceptional cases. However, our implementation has not been optimized for speed and memory consumption. In general, DeNovoΔPos is more time-consuming. Half of the spectra were analyzed within about 2 min. The running time of our algorithm depends on the size of the spectrum graph. The algorithms considered two masses to be equal if they differ by at most 0.02 Da. Moreover, a simple merging algorithm is applied during the construction of the spectrum graph to reduce the size of the graph as described in [14]. We observed a great variation of spectrum graph sizes in our experiments. The spectrum graphs contained roughly 8400 edges on average, whereas the largest observed graph contained 23,000 edges. Spectra measured on low resolution lead to denser spectrum graph, i.e. to a larger number of edges, but a lower number of vertices. However, we did not study the performance and runtime of our algorithms on this type of spectra.

## Conclusion

In this paper, we propose the first algorithms for exploiting the retention time information in de novo peptide sequencing. We study three retention time prediction models and develop algorithms for computing a sequence that matches the experimental mass spectrum as well as possible and is in accordance with the observed retention time. The experimental evaluation of our algorithms shows that identification rates can definitely be improved by exploiting this additional information. Yet, the proposed algorithms score sequences with a very simplistic scoring function that only counts explained and measured masses and does not consider any other available information. For real-world applications, a more evolved scoring function using all available information needs to be integrated. While [14] introduces a new scoring model, we explore ways of exploiting the retention time information. The proposed algorithms open room for developing new scoring functions that consider both the retention time information and the symmetric difference scoring model.

## Authors' contributions

Preliminary work of this study was carried out by YV and VV during their bachelor theses. All authors contributed to the design of the algorithms and

the experiments. YV, VV, and TT implemented the software. TT performed the experiments. All work was guided by TH in the whole process. TT wrote most parts of the manuscript with the help of TH. All authors contributed to the writing of the paper. All authors read and approved the final manuscript.

#### Acknowledgements

We would like to thank Peter Widmayer, Christian Panse, Witold Wolski, and Ludovic Gillet for helpful discussions. Moreover, we thank the anonymous reviewers of the preliminary conference version of this study for their constructive criticism.

#### Competing interests

The authors declare that they have no competing interests.

#### Availability of data and materials

The datasets analyzed during the current study are available in the peptide atlas repository, <http://www.peptideatlas.org> with dataset ID PASS00289 [15]. The implementation of our algorithm is available under a BSD license at <http://www.github.com/tschager>.

#### Consent for publication

Not applicable.

#### Ethics approval and consent to participate

Not applicable.

#### Funding

Funding information is not applicable.

#### Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Received: 30 October 2017 Accepted: 20 August 2018

Published online: 29 August 2018

#### References

- Kinter M, Sherman NE. Protein sequencing and identification using tandem mass spectrometry. New York: Wiley; 2000. <https://doi.org/10.1002/0471721980>.
- Hughes C, Ma B, Lajoie GA. De novo sequencing methods in proteomics. *Proteome Bioinform*. 2010;6(4):105–21. [https://doi.org/10.1007/978-1-60761-444-9\\_8](https://doi.org/10.1007/978-1-60761-444-9_8).
- Van Riper SK, de Jong EP, Carlis JV, Griffin TJ. Mass spectrometry-based proteomics: basic principles and emerging technologies and directions. *Adv Exp Med Biol*. 2013;990:1–35. [https://doi.org/10.1007/978-94-007-5896-4\\_1](https://doi.org/10.1007/978-94-007-5896-4_1).
- Dančik V, Addona TA, Clauser KR, Vath JE, Pevzner PA. De novo peptide sequencing via tandem mass spectrometry. *J Comput Biol*. 1999;6(3–4):327–42. <https://doi.org/10.1089/106652799318300>.
- Chen T, Kao M-Y, Tepel M, Rush J, Church GM. A dynamic programming approach to de novo peptide sequencing via tandem mass spectrometry. *J Comput Biol*. 2001;8(3):325–37. <https://doi.org/10.1089/10665270152530872>.
- Jeong K, Kim S, Pevzner PA. Uninovo: a universal tool for de novo peptide sequencing. *Bioinformatics*. 2013;29(16):1953–62.
- Ma B, Zhang K, Hendrie C, Liang C, Li M, Doherty-Kirby A, Lajoie G. Peaks: powerful software for peptide de novo sequencing by tandem mass spectrometry. *Rapid Commun Mass Spectrom*. 2003;17(20):2337–42. <https://doi.org/10.1002/rcm.1196>.
- Ma B. Novor: Real-time peptide de novo sequencing software. *J Am Soc Mass Spectrom*. 2015;26(11):1885–94. <https://doi.org/10.1007/s13361-015-1204-0>.
- Shinoda K, Sugimoto M, Tomita M, Ishihama Y. Informatics for peptide retention properties in proteomic LC–MS. *Proteomics*. 2008;8(4):787–98. <https://doi.org/10.1002/pmic.200700692>.
- Moruz L, Käll L. Peptide retention time prediction. *Mass Spectrom Rev*. 2016;36(5):615–23. <https://doi.org/10.1002/mas.21488>.
- Palmlad M, Ramstöm M, Markides KE, Håkansson P, Bergquist J. Prediction of chromatographic retention and protein identification in liquid chromatography/mass spectrometry. *Anal Chem*. 2002;74(22):5826–30. <https://doi.org/10.1021/ac0256890>.
- Strittmatter EF, Kangas LJ, Petritis K, Mottaz HM, Anderson GA, Shen Y, Jacobs JM, Camp DG, Smith RD. Application of peptide LC retention time information in a discriminant function for peptide identification by tandem mass spectrometry. *J Proteome Res*. 2004;3(4):760–9. <https://doi.org/10.1021/pr049965y>.
- Frank Y, Hruz T, Tschager T, Venzin V. Improved de novo peptide sequencing using LC retention time information. In: Schwartz R, Reinert K, eds. 17th international workshop on algorithms in bioinformatics (WABI 2017), vol. 88. Leibniz International Proceedings in Informatics (LIPIcs) Dagstuhl, Germany: Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik; 2017. p. 26–617.
- Tschager T, Rösch S, Gillet L, Widmayer P. A better scoring model for de novo peptide sequencing: the symmetric difference between explained and measured masses. *Algorithms Mol Biol*. 2017. <https://doi.org/10.1186/s13015-017-0104-1>.
- Röst HL, Rosenberger G, Navarro P, Gillet L, Miladinović SM, Schubert OT, Wolski W, Collins BC, Malmström J, Malmström L, Aebersold R. Open-swath enables automated, targeted analysis of data-independent acquisition ms data. *Nat Biotechnol*. 2014;32(3):219–23. <https://doi.org/10.1038/nbt.2841>.
- Röst HL, Sachsenberg T, Aiche S, Bielow C, Weisser H, Aicheler F, Andreotti S, Ehrlich H-C, Gutenbrunner P, Kenar E, Liang X, Nahnsen S, Nilse L, Pfeuffer J, Rosenberger G, Rurik M, Schmitt U, Veit J, Walzer M, Wojnar D, Wolski WE, Schilling O, Choudhary JS, Malmström L, Aebersold R, Reinert K, Kohlbacher O. Openms: a flexible open-source software platform for mass spectrometry data analysis. *Nat Methods*. 2016;13(9):741–8. <https://doi.org/10.1038/nmeth.3959>.
- Guo D, Mant CT, Taneja AK, Parker JMR, Rodges RS. Prediction of peptide retention times in reversed-phase high-performance liquid chromatography I. Determination of retention coefficients of amino acid residues of model synthetic peptides. *J Chromatogr A*. 1986;359:499–518. [https://doi.org/10.1016/0021-9673\(86\)80102-9](https://doi.org/10.1016/0021-9673(86)80102-9).
- Krokhin OV, Craig R, Spicer V, Ens W, Standing KG, Beavis RC, Wilkins JA. An improved model for prediction of retention times of tryptic peptides in ion pair reversed-phase HPLC: its application to protein peptide mapping by off-line HPLC-MALDI MS. *MCP*. 2004;3(9):908–19. <https://doi.org/10.1074/mcp.M400031-MCP200>.
- Krokhin OV. Sequence-specific retention calculator. Algorithm for peptide retention prediction in ion-pair RP-HPLC: application to 300- and 100-Å pore size c18 sorbents. *Anal Chem*. 2006;78(22):7785–95. <https://doi.org/10.1021/ac060777w>.
- Spicer V, Grigoryan M, Gotfrid A, Standing KG, Krokhin OV. Predicting retention time shifts associated with variation of the gradient slope in peptide RP-HPLC. *Anal Chem*. 2010;82(23):9678–85. <https://doi.org/10.1021/ac102228a>.
- Schubert OT, Gillet LC, Collins BC, Navarro P, Rosenberger G, Wolski WE, Lam H, Amodei D, Mallick P, MacLean B, Aebersold R. Building high-quality assay libraries for targeted analysis of swath MS data. *Nat Protocols*. 2015;10(3):426–41. <https://doi.org/10.1038/nprot.2015.015>.
- Eng JK, Jahan TA, Hoopmann MR. Comet: an open-source ms/ms sequence database search tool. *Proteomics*. 2013;13(1):22–4. <https://doi.org/10.1002/pmic.201200439>.
- Keller A, Nesvizhskii AI, Kolker E, Aebersold R. Empirical statistical model to estimate the accuracy of peptide identifications made by ms/ms and database search. *Anal Chem*. 2002;74(20):5383–92. <https://doi.org/10.1021/ac025747h>.
- R Core Team. R: a language and environment for statistical computing. R foundation for statistical computing, Vienna, Austria; 2017. <https://www.R-project.org/>