



**Fragmentation-free LC-MS can identify hundreds of proteins.**

Journal:	<i>PROTEOMICS</i>
Manuscript ID:	pmic.200900765.R1
Wiley - Manuscript type:	Research Article
Date Submitted by the Author:	10-Jul-2010
Complete List of Authors:	Bochet, Pascal; Institut Pasteur, System Biology; CNRS, U2171 Rügheimer, Frank; Institut Pasteur, System Biology Guina, Tina; Univ. of Washington, Dept of Pediatrics Brooks, Peter; CNRS, UMR7592; Université Denis Diderot, Institut Jacques Monod Goodlett, David; Univ. of Washington, Dept of Medicinal Chemistry Clote, Peter; Boston College, Biology Dept; Ecole Polytechnique, LIX Schwikowski, Benno; Institut Pasteur, Laboratory for System Biology, Department Genomes and Genetics; CNRS, URA2171
Key Words:	Protein automated identification, Mass spectrometry, Peptide mass fingerprinting, Liquid chromatography-tandem mass spectrometry



# Fragmentation-free LC-MS can identify hundreds of proteins.

Pascal Bochet<sup>1,2</sup>, Frank Rügheimer<sup>1,2</sup>, Tina Guina<sup>3</sup>, Peter Brooks<sup>4,5</sup>,  
David Goodlett<sup>6</sup>, Peter Clote<sup>7,8,9</sup>, Benno Schwikowski<sup>1,2</sup>.

1 Institut Pasteur, Laboratoire de Biologie Systémique, Dept Génomes et Génétique, Paris, France

2 CNRS, URA2171, Paris, France

3 University of Washington, Dept of Pediatrics, Seattle, WA, USA

4 CNRS, UMR7592, Institut Jacques Monod, Paris, France

5 Université Paris 7, Paris, France

6 University of Washington, Dept of Medicinal Chemistry, Seattle, WA, USA

7 Biology Dept, Boston College, Boston, MA, USA

8 Ecole Polytechnique, LIX, Palaiseau, France

9 LRI, Université Paris XI, Orsay, France

Authors for correspondence: Pascal Bochet, Benno Schwikowski

Address: Laboratoire de Biologie Systémique, Institut Pasteur, 25, rue du Docteur Roux, 75015 Paris, France

email: pascal.bochet@pasteur.fr

Fax: 33 (0)1.40.61.37.04

## Keywords

Automated protein identification, Mass spectrometry, Peptide mass fingerprinting, Liquid chromatography-tandem mass spectrometry.

Length: 8122 words

## Abstract

One of the most common approaches for large-scale protein identification is liquid chromatography, followed by mass spectrometry (LC-MS). If more than a few proteins are to be identified, the additional fragmentation of individual peptides has so far been considered as indispensable, and thus, the associated costs, in terms of instrument time and infrastructure, as unavoidable. Here, we present evidence to the contrary. Using a combination of [1] highly accurate and precise mass measurements,[2] modern retention time prediction, and [3] a robust scoring algorithm, we were able to identify 257 proteins of *Francisella tularensis* from a single LC-MS experiment in a fragmentation-free approach (i.e., without experimental fragmentation spectra). This number amounts to 59% of the number of proteins identified in a standard fragmentation-based approach, when executed with the same false discovery rate. Independent evidence supports at least 27 of a set of 31 proteins that were identified only in the fragmentation-free approach. Our results suggest that additional developments in retention time prediction, measurement technology, and scoring algorithms may render fragmentation-free approaches an interesting complement or an alternative to fragmentation-based approaches.

# 1 Introduction

## 1.1 Protein identification by HPLC-MS

One of the key goals of proteomics experiments is the identification of proteins from samples of unknown composition. Thanks to continuous efforts in genome sequencing, genome and protein sequence information are now available for an increasing number of species [1]. For those species, the identification of proteins is accomplished by matching the results of Mass Spectrometry (MS) measurements against a sequence database of candidate proteins.

When only few protein are in the sample, a well-established procedure, Peptide Mass Fingerprinting (PMF), can be applied for protein identification [2, 3, 4, 5]. The protein is digested by an enzyme with known substrate specificity, the masses of the resulting peptides are measured by MS and the relevant sequence database is searched for proteins which generate a corresponding set of masses after enzymatic cleavage. Unfortunately, for a larger number of different proteins in a sample, the combinatorial explosion of possible mass combinations prevents reliable identification by PMF.

Additional information can be obtained by further fragmenting peptide ions isolated from the first MS analysis. The fragment masses are commonly determined in second MS steps interleaved with the first MS. The series of observed fragments is compared to predictions made from the known sequences and the most likely peptide candidates are selected [6]. One problem with this strategy is the high number of peptide ions generated by the first MS, which precludes their exhaustive analysis by fragmentation. This is partly overcome by separating the peptides prior to MS, usually on a HPLC column [7]. Most often, Reverse Phase columns are used, which separate the peptides mainly on the basis of hydrophobic interactions with the column [8]. Although HPLC separation helps alleviate the problem, not all peptides can be

1  
2  
3  
4  
5  
6  
7  
8 analyzed and identified by fragmentation.

9  
10 Apart from spreading peptides over the time domain, the preliminary HPLC  
11 analysis also provides for each of them a retention time (RT). This readily available,  
12 additional information has so far rarely been exploited for identification. Here we  
13 demonstrate that, together with highly accurate peptide mass measurements, the  
14 additional retention time information enables the identification of hundreds of  
15 proteins even without the use of MS/MS fragmentation.  
16  
17  
18  
19  
20  
21  
22  
23

## 24 **1.2 Using retention time and accurate peptide** 25 **mass for identification** 26 27 28 29 30

31 Several previous efforts have used HPLC and peptide retention time for the  
32 identification of proteins. We will briefly discuss some of them.  
33  
34  
35  
36  
37

### 38 **1.2.1 Theoretical justification** 39 40

41 A recent theoretical study by Norbeck *et al.* [9] has given insight into the precision of  
42 mass and retention time determination required for fragmentation-free identification  
43 of tryptic peptides. The authors evaluated to what extent the determination of mass  
44 and HPLC retention time allows unique identification of tryptic peptides. In  
45 particular, they studied the number of uniquely identifiable peptides as a function of  
46 the precision of mass and retention time and the size of the relevant proteome. The  
47 theoretical study showed that the addition of the HPLC retention time strongly  
48 enhances the identification power of MS. For instance, for a peptide mass of around  
49 2250 Da, in the bacterium *Deinococcus radiodurans* (3167 proteins), approximately  
50 20% of the peptides could be uniquely identified at a mass error rate of 5 parts per  
51 million (ppm). Combining the mass measurement with the HPLC retention time  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 (determined with a 1% error rate) increases the proportion of unique peptides from  
9  
10 20% to 90%. For a more complex proteome (*Homo sapiens*) the proportion of unique  
11  
12 peptides would be 2% at 5 ppm, and would increase to 50% with the additional use of  
13  
14 retention time prediction with 1% error. The study of Norbeck *et al.* demonstrates that  
15  
16 taking into account the retention time of the peptides when searching for matches in  
17  
18 the database can drastically reduce the number of candidate peptides. It does not,  
19  
20 however, address practical implementation aspects of this idea.  
21

### 22 23 24 **1.2.2 Determining peptide retention time** 25

26  
27 HPLC peptide retention times can be used for peptide identification by comparing  
28  
29 experimentally obtained retention times with retention times that have been predicted,  
30  
31 or previously measured.  
32

### 33 34 **Measuring retention time** 35

36  
37 A straightforward approach to obtain the retention times of a small number of  
38  
39 peptides in a given HPLC separation system is to measure them: Synthetic peptides  
40  
41 with the desired sequences are injected and detected when they exit from the column  
42  
43 (see for instance [10]). However, this approach is limited to peptides that can be  
44  
45 synthesized or highly purified and can thus not provide the hundreds or thousands of  
46  
47 retention times required for protein identification from complex mixtures. With the  
48  
49 advent of LC-MS/MS methods, the retention times of peptides selected for analysis  
50  
51 by fragmentation after HPLC became available in very large amounts. These  
52  
53 retention times can then be used directly or serve as training examples for diverse  
54  
55 retention time prediction algorithms. But to account for unavoidable experimental  
56  
57 variability, the retention times need to be normalized. For this purpose Strittmatter *et*  
58  
59 *al.* [11], for instance, obtained data from individual runs and normalized the retention  
60  
times by affine transformations optimized over the whole dataset. Still, the intrinsic

1  
2  
3  
4  
5  
6  
7  
8 variability of HPLC places a limit on the precision that can be obtained for retention  
9  
10 time.

### 11 **Predicting retention time**

12  
13  
14  
15 An alternative to measuring peptide retention times is to predict them. Methods to  
16  
17 predict peptide retention times from their primary sequences have been extensively  
18  
19 reviewed recently by Baczek and Kaliszan [8].

20  
21 One early approach was based on data from few injected synthetic peptides [10].  
22  
23 This work was limited by the amount of data available at the time, but introduced the  
24  
25 idea that the hydrophobic interaction of a peptide with the matrix of the column could  
26  
27 be predicted from the amino-acid composition by a linear model.

28  
29 Later, other work relied on the idea of such a linear model, either using data from  
30  
31 a rather small dataset of synthetic peptides [12] or on naturally occurring tryptic  
32  
33 peptides [13]. Other peptides properties (e.g.its length [14]) have also been used with  
34  
35 the linear model.

36  
37 Other approaches have used artificial neural networks [15, 16] on a very large  
38  
39 number of peptides (345000 peptides) or support vector machines [17, 18] on a much  
40  
41 smaller dataset to obtain predictive models.

42  
43 Krokhin's SSRCalc tool [19], develops the idea of a linear model further by  
44  
45 employing different hydrophobicity coefficients for the amino acids close to the  
46  
47 termini of the peptide, and by taking into account the length of the peptide, its likely  
48  
49 secondary structure and effects between neighboring aminoacids. Contrary to many  
50  
51 other models, SSRCalc does not require experiment or platform-specific training of  
52  
53 the model. For each peptide, SSRCalc provides a hydrophobicity index which can be  
54  
55 converted into a predicted retention time by an affine transformation.

### 56 **1.2.3 Using retention time for peptide/protein identification**

57  
58  
59  
60 Regardless of how retention times are determined, several strategies have been

1  
2  
3  
4  
5  
6  
7  
8 developed for their use in peptide and protein identification.

9  
10 The Accurate Mass and Time tag (AMT) strategy, introduced by the group of  
11 R.Smith [11], uses the retention times of the peptides recorded in MS/MS  
12 experiments. Simultaneous matching of the mass and the retention times are then  
13 used to identify the peptides in a specific database and, subsequently, the proteins.  
14 The advantages of this method are high proteome coverage (at least for small  
15 proteomes like *Deinococcus radiodurans*) and short analysis time (3 hours, i.e. the  
16 typical duration of a HPLC-MS run). But this strategy requires that the retention  
17 times of thousands of peptides are recorded before the identification phase.  
18 Strittmatter *et al.* [11], for instance, recorded the retention times of more than 12000  
19 peptides from 1067 runs of MS/MS experiments.  
20  
21  
22  
23  
24  
25  
26  
27  
28

29 Instead of using previously measured retention times, only available for peptides  
30 characterized in previous experiments, several authors have used retention times  
31 predicted from the primary peptide sequence. Retention time predictions are generally  
32 available for every peptide, but the accuracy of this prediction will naturally be  
33 limited by the underlying model.  
34  
35  
36  
37  
38

39 Palmblad *et al.* [13] used their composition-based, linear, retention time predictor.  
40 The deviation from predicted mass and retention time is combined into a Euclidean  
41 distance after separate normalization in each dimension. The quality of the peptide  
42 match and the identification is estimated from the value of this distance. In later work  
43 the mass and retention time determinations were combined using a likelihood-based  
44 method [20].  
45  
46  
47  
48  
49  
50

51 Similarly, using their Support Vector Machine predictor, Pfeifer *et al.* [17, 18]  
52 added to the fragmentation-based detection a filter based on the deviation of the  
53 predicted retention times from the observed and normalized values. In this approach,  
54 a potential peptide match detected on the basis of fragmentation spectra is excluded if  
55 the deviation between the observed and predicted peptide retention times is above a  
56  
57  
58  
59  
60



1  
2  
3  
4  
5  
6  
7  
8 given threshold. The threshold value itself is determined via a parametrized linear  
9  
10 model that allows the acceptable absolute deviation to increase with normalized  
11  
12 observed retention time.

13  
14 Krokhin *et al.* [19] used their variant of the linear predictor in two different  
15  
16 methods. MART analysis is based on fragmentation data from HPLC-MS/MS and  
17  
18 uses predicted retention time to progressively eliminate peaks from the spectra, and  
19  
20 thus speed up the computational identification process [21]. A second method is based  
21  
22 only on HPLC-MS data. It first identifies the proteins compatible with all peptide  
23  
24 masses detected in the MS data using the PMF search program Profound [5]. It then  
25  
26 uses the extent of the correlation between predicted and measured retention times of  
27  
28 the peptides belonging to each single protein to filter the results of PMF protein  
29  
30 identification (see [http://hs2.proteome.ca/SSRCalc/OKrokhin\\_ASMS07.pdf](http://hs2.proteome.ca/SSRCalc/OKrokhin_ASMS07.pdf)).  
31  
32  
33

### 34 **1.3 A novel fragmentation-free approach**

35  
36  
37  
38 Except for Krokhin's correlation method, all identification strategies discussed above  
39  
40 depend upon previously measured fragmentation spectra, which are required either to  
41  
42 directly establish peptide-specific retention times, or to train predictive models. Thus,  
43  
44 one way to improve these methods is to rely on more reproducible HPLC or more  
45  
46 accurate predictions. However the reproducibility of HPLC and the precision of  
47  
48 predictions are limited, depend on the particular experimental setup, and cannot be  
49  
50 improved beyond the intrinsic variability of the HPLC procedure.

51  
52 Therefore, we designed a new approach which reduces the dependency on the  
53  
54 retention time precision. Our approach relies on the retention *order* of the peptides  
55  
56 rather than their retention time. By substituting the retention times by the retention  
57  
58 order, we extract a more robust characterization of the retention process. Furthermore,  
59  
60 we do not seek to identify isolated peptides, but rather groups of peptides coming

1  
2  
3  
4  
5  
6  
7  
8 from the digest of the same protein: Each peptide will be identified from the  
9  
10 experimental data only if other predicted peptides from the same candidate protein are  
11  
12 also found in their predicted order. Unlike most of the methods described above, the  
13  
14 proposed approach identifies proteins in one single step and not by first finding  
15  
16 peptides and then proteins.

17  
18 These ideas have been implemented in an alignment method (detailed in Methods  
19  
20 **2.3.2**) at the core of our identification strategy based on what we have called the  
21  
22 Ordered Peptide Match score (OPM score).  
23  
24

## 25 **2 Material and methods**

### 26 27 28 29 **2.1 Principle**

30  
31  
32  
33 The principle of the new method is illustrated by Figure 1. For each protein in the  
34  
35 relevant sequence database we first predict the peptides resulting from its enzymatic  
36  
37 digest. In a second step, the HPLC retention time predictions for the peptides are  
38  
39 obtained and the peptides are sorted by predicted retention time. In a third step we  
40  
41 attempt to assign each predicted peptide for a given protein to a peak in a spectrum in  
42  
43 the predicted order. We solved this alignment problem by dynamic programming.  
44  
45  
46  
47  
48

49  
50 Figure 1 near here  
51  
52  
53

### 54 **2.2 Experimental data**

55  
56  
57  
58 Experimental data were obtained from the species *Francisella tularensis subsp.*  
59  
60 *novicida*, a Gram negative bacterium that is a causative agent of the human and

1  
2  
3  
4  
5  
6  
7  
8 animal disease tularemia<sup>1</sup>.  
9

### 10 **2.2.1 Sample preparation**

11  
12 Bacteria were grown in rich medium, harvested by centrifugation during the  
13 exponential growth phase and were broken with ultrasound after several rounds of  
14 freezing in dry ice-ethanol and thawing in a water bath at 12°C. Membrane-bound  
15 and soluble proteins were separated by centrifugation. Aliquots of both preparations  
16 were solubilised in 6M urea, reduced with tris(2-carboxyethyl)phosphine (TCEP) and  
17 alkylated with iodoacetamide. Samples were then treated with DTT and digested with  
18 trypsin. Samples were desalted on C18 reverse phase columns and analyzed by  
19 HPLC-MS/MS. For more details see [23].  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31

### 32 **2.2.2 Sample analysis**

33  
34 Samples were loaded on a C-18 200Å pore size precolumn and separated on a C-18  
35 100Å pore size microcolumn (75 µm × 11 cm, 5 µm)(Magic C-18 AQ beads,  
36 Michrom Bioresources, CA). Retention was achieved by a 60 min long linear gradient  
37 from 10% to 45% of solvent B (100% acetonitrile) in solvent A (0.1% formic acid,  
38 5% acetonitrile) at a flow rate of 200nl/min. The spectra were obtained in a LTQ-FT-  
39 ICR-MS (Linear Trap Quadrupole-Fourier Transform-Ion Cyclotron Resonance Mass  
40 Spectrometer, Thermo Electron, San Jose, CA). Ions with m/z between 400 and 1800  
41 were analyzed by FT-ICR. The most intense ions were fragmented by CID (Collision  
42 Induced Dissociation) and analyzed by LTQ. The nucleotide and protein sequences  
43 for *Francisella tularensis subsp. novicida* are publicly available from NCBI under the  
44 accession number NC\_008601.1.  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

---

<sup>1</sup> *Francisella tularensis* is considered a potential bioterrorism agent, see [22].

## 2.3 Computational methods

### 2.3.1 Data reduction

In a first data analysis step, experimental data were preprocessed. The aim of this step was to reduce the spectra to mono-isotopic, mono-charged spectra. A method adapted from Matthiesen for his software VEMS [24] was applied, with the main difference that each spectrum was treated independently from the others.

Deisotoping was carried out as follows: a global maximum charge was chosen for the ions. The spectrum was searched for peaks with intensities above a set threshold, starting with low values of the *mass/charge* ratio. When found, such peaks were considered as putative mono-isotopic peaks and a series of matching isotopic peaks with *m/z* spacing of  $1/(\text{charge} \pm \text{tolerance})$  was searched. If found, such a series was removed from the spectrum, stored with the corresponding charge value, and the search was continued. When the higher end of the spectrum was reached, the process was repeated from the lower end with the next lower value of the ion charge. The resulting list of mono-isotopic, mono-charged peaks was stored and used for searching the predicted peptides from the protein digest.

To reduce the effect of experimental noise we only considered peaks with intensities in the top 24.5% quantile of each spectrum. In our experimental data this threshold was found to remove low intensity noise efficiently.

### 2.3.2 Computation of Ordered Peptide Match score (OPM)

The complete *in silico* digest was performed on all proteins predicted by the database. Only peptides at least four amino acids long were kept. For each protein these peptides were sorted according to their predicted order of retention from the HPLC column. For the results reported here we have used the hydrophobic predictions in TFA of the SSRCalc software, version 3.2, available on-line at

1  
2  
3  
4  
5  
6  
7  
8 <http://hs2.proteome.ca/SSRCalc/SSRCalc32.html> [21, 25] .  
9

10 For each protein, the algorithm described in Figure 2 computes the optimal  
11 alignment of the predicted peptides with the ions found in the experimental spectra.  
12 This maximizes the OPM score. The score optimized in the alignment is defined as  
13 follows: the mass of each predicted peptide is compared to the closest mass in each  
14 spectrum after deisotoping and decharging. If  $\Delta$ , the difference between the two  
15 masses is less than a given threshold, the contribution of the match to the score is +1,  
16 otherwise a penalty is counted (-1). The alignment is then built in a progressive way,  
17 as shown by Figure 2. A value of 0.008 Da for the mass threshold has been chosen  
18 because 95% of the peptides detected by fragmentation are within this distance from  
19 the true mass (see below **2.4**).  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30

31 For the particular choice of scoring function used here, the optimal OPM score is  
32 an affine transform of the number of predicted peptides matched to the experimental  
33 data. Therefore in the rest of the text we use this property to simplify the analysis and  
34 the figures by using the number of matching peptides.  
35  
36  
37  
38

39 Figure 2 near here  
40  
41  
42  
43

### 44 **2.3.3 Interpretation of the Ordered Peptide Match scores**

45

46 To evaluate the statistical significance of the OPM scores obtained for the proteins  
47 from the sequence databases, they were compared to the OPM scores of decoy  
48 proteins obtained by random permutation of the sequences of the real ones. Each of  
49 these decoy proteins had the same length, number of tryptic peptides and amino acid  
50 composition as its naturally occurring counterpart. In the database the number of  
51 proteins with the largest content of digest peptides was low. To obtain a sufficiently  
52 accurate estimate of the quantile limits for these large proteins (with more than 56  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 peptides), several permutations were performed for each of them.

9  
10 Because the alignment scores depend on the number of tryptic peptides (see  
11 Section **3.1** and Figure 3), and in the absence of a substantial theoretical  
12 justification for the use of a particular parametric distribution model, we decided to  
13 apply quantile regression [26, 27] to describe the empirical distribution of scores for  
14 the decoy proteins. Briefly, for each value of the variable (here, the number of tryptic  
15 peptides in the protein) quantile regression uses a weighting function to construct a  
16 local score distribution for the randomized proteins. The score of each real protein  
17 was compared to that local distribution for the same number of tryptic peptides, and  
18 the corresponding quantile was reported as the significance of the match.  
19  
20  
21  
22  
23  
24  
25  
26  
27

#### 28 29 **2.3.4 Estimation of the false discovery rate**

30  
31  
32 The observed distribution of scores for real proteins is modeled as a mixture  
33 generated by drawing from two populations: the proteins absent from the sample and  
34 those present. Since neither the relative contributions of each group nor the actual  
35 distribution of scores for proteins in the sample are known, we applied the following  
36 approximation. For each quantile limit, we considered that all proteins with scores  
37 below the limit were absent from the sample. We assumed a score distribution for  
38 these absent proteins identical to that of the decoy proteins, and that no real protein  
39 effectively present in the sample had a score below the quantile limit. This amounts to  
40 a conservative choice of the largest possible number of absent proteins for a given  
41 total number of proteins. Based on this distribution we calculated the number of  
42 absent real proteins expected to have a score above the quantile limit (false positives).  
43  
44 For each quantile limit this allowed us to estimate the false discovery rate (FDR) as  
45 the proportion of false positives among the total positives with scores above the limit.  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 2.4 Comparison with fragmentation

As the experimental data we used also included fragmentation spectra, we were able to compare the set of proteins identified by the OPM method with the set of proteins obtained by the standard fragmentation-based approach on the same experiments. Note that this analysis is likely to slightly underestimate the power of a pure OPM approach: If no time were devoted to the acquisition of fragmentation spectra, the additionally acquired MS spectra could be expected to improve the performance of the OPM method. Since this effect may not be very large, and since the comparison on the same dataset avoids statistical differences in all other respects, we chose to use the same dataset to evaluate both approaches. Fragmentation spectra were analyzed with the Mascot software [4]. Version 2.2.0 was used with the following parameters: sequence database: *Francisella* NC\_008601 (NCBI), mass type: monoisotopic, fixed modifications: Carbamidomethyl (C), no variable modifications, enzyme: Trypsin, no missed cleavages, peak mass tolerance: 0.6, precursor mass tolerance: 1.2. For comparability, the false discovery rate (FDR) of Mascot was estimated by using a decoy database as described by [28] and the number of positive proteins retained from Mascot was set at 436 to obtain a FDR equal to that of the OPM method (5.7%).

## 3 Results

### 3.1 Correlation of peptide number and score

Figure 3 shows, for each *Francisella* protein and each decoy protein, the OPM scores as a function of the number of their tryptic peptides (data in Table 1). The basis of the analysis is the comparison between the scores of the proteins encoded by the genome of the bacteria (1719 real proteins) and the decoy database made of 2205 fictive

1  
2  
3  
4  
5  
6  
7  
8 proteins.

9  
10 For both groups of proteins, the score is, in a first approximation, proportional to  
11 the number of tryptic peptides in the protein. However, the slope of the associated  
12 least-squares regression line is higher for the real proteins (slope:  $0.164 \pm 0.003$ ) than  
13 for the randomized ones ( $0.116 \pm 0.002$ ). This reflects the fact that the number of  
14 peptides matching with the experimental data is higher for the real proteins than for  
15 the decoy proteins. Also noteworthy is the larger deviation from linearity of scores for  
16 real proteins ( $R^2 = 0.567$ ) than for the decoy proteins ( $R^2 = 0.673$ ). This reflects the  
17 presence of real proteins with a large number of matching peptides and a high score  
18 absent from the decoy database.  
19  
20  
21  
22  
23  
24  
25  
26

27 Figure 3 near here  
28  
29  
30  
31  
32  
33

### 3.2 Detection of proteins and false discovery rate

34  
35

36 In this section we make use of the OPM score described above to classify proteins as  
37 either present or absent from the experimental sample. Because the alignment scores  
38 depend on the number of peptides in the proteins we could not use it directly for this  
39 purpose. Therefore, we used a quantile regression method. By building a conditional,  
40 empirical, local distribution of the scores for each number of peptides in the proteins,  
41 quantile regression allows the comparison of scores for real proteins with those of  
42 corresponding randomized proteins. The proteins from the experiment shown in  
43 Figure 3 were classified as described in methods (section **2.3.3** and see quantile  
44 assignment in Table 1). For several quantile values, the number of proteins with score  
45 above the quantile threshold is indicated in Table 2, top. Table 2 shows that the  
46 number of detected proteins (real proteins in the table) is much larger than the  
47 corresponding number from randomized proteins (random proteins in the table). The  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3  
4  
5  
6  
7  
8 proportion of randomized proteins above each threshold (random proteins) is  
9 consistent with the corresponding quantile limit, thus showing the absence of any  
10 large systematic bias in the local quantile estimator.  
11  
12

13 Estimates of the FDR, the proportion of false positives among the positives, are  
14 also shown in Table 2. This estimated proportion decreases with the stringency of the  
15 detection from 34% for a quantile of 80% (721 proteins detected) to 5.7% for 99%  
16 (257 proteins detected).  
17  
18  
19  
20  
21  
22  
23

24 Table 2 near here  
25  
26  
27  
28

### 29 **3.3 Comparison of ordered and unordered** 30 **peptide scores.** 31 32 33

34  
35  
36 When compared to classical peptides mass fingerprinting (PMF), the OPM scoring  
37 approach introduced here additionally considers the degree of consistency between  
38 the orders of measured and predicted retention times. The rationale for the present  
39 work is the idea that the order of elution of peptides based on their predicted retention  
40 times will introduce additional discriminatory power for matching the peptides. To  
41 verify the impact of this additional element, we performed the same analysis as  
42 above, but without taking into account the elution order of the peptides. Specifically,  
43 we created a fictive data set, where all experimental peaks retained after decharging  
44 and deisotoping (cf. Section 2.3.1) were present in every spectrum. The alignment  
45 procedure was then performed against this new data set, effectively neutralizing the  
46 order information present in the original data.  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56

57  
58 As with the ordered peptides, the scores are proportional to the number of tryptic  
59 peptides. As expected, the alignment scores obtained (unordered peptide match  
60

1  
2  
3  
4  
5  
6  
7  
8 (UPM) scores) were larger than the OPM scores computed before. They were higher  
9  
10 for the real proteins (slope:  $0.427 \pm 0.005$ ) than for the randomized ones (slope:  $0.356$   
11  $\pm 0.003$ ). The deviation from proportionality is less than that for the ordered peptides  
12 ( $R^2 = 0.807$  for real proteins,  $R^2 = 0.858$  for randomized). In summary, the score  
13 difference between the real and the randomized proteins, which is the basis for the  
14 protein recognition procedure introduced here, is much weaker without the order than  
15 with it.  
16  
17  
18  
19  
20  
21

22 This is also reflected in the detection of proteins as shown at the bottom of Table  
23  
24 2. In this table we can see that for every value of the quantile limit, the number of real  
25 proteins labelled as present in the sample is lower than that obtained with the ordered  
26 peptides match score. In addition, the estimated FDR is larger. This is further  
27 illustrated in Figure 4 A, where the FDR is plotted against the number of detected  
28 proteins: For almost all numbers of recognized proteins (obtained by varying the  
29 quantile limit) the estimated FDR is substantially smaller when based on the OPM  
30 score than the UPM score.  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41

42 Figure 4 near here  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53

### 54 **3.4 Comparison of ordered peptide matched score** 55 **with fragmentation based identification.** 56 57 58 59 60

We determined how many of the proteins identified using the OPM method based on retention time and mass were also detected by Mascot on the basis of the fragmentation and MS/MS data from the same experiment. With our Mascot settings (see Methods), 436 proteins were identified in the experiment, out of 1719 annotated

1  
2  
3  
4  
5  
6  
7  
8 in the genome of *F. tularensis* (see Table 3). Figure 4B shows the ROC curves  
9  
10 obtained by comparing the identification results based on ordered or unordered  
11  
12 peptide match scores to the proteins identified by Mascot. It can be seen that order  
13  
14 information consistently improves the performance of the protein detector for values  
15  
16 of the false positive rate (FPR) up to FPR = 0.4 (corresponding to a True Positive  
17  
18 Rate of 0.91, at the low quantile limit of 0.66 (OPM), 0.654 (UPM)). This  
19  
20 improvement is coupled to a significant reduction in the estimated false discovery rate  
21  
22 (cf Figure 4A.). Table 2, Top also shows that the OPM method can detect proteins not  
23  
24 identified from the fragmentation data: At the 99% quantile limit, 257 proteins were  
25  
26 identified, of which 226 had been also detected by fragmentation and Mascot,  
27  
28 revealing 31 newly detected proteins. Even when peptides of these 31 proteins had  
29  
30 been identified they did not support the identification of the protein by Mascot under  
31  
32 the settings used. For 6 of the 31 newly found proteins, no peptide had been identified  
33  
34 by Mascot. For 8 of them, only one peptide had been identified by Mascot, and, for  
35  
36 the remaining 17 proteins, at least 2 and up to 10 peptides had been identified. In  
37  
38 addition, in a more extensive characterization of proteins expressed in *Francisella*  
39  
40 *tularensis* (1080 proteins detected), 27 of these 31 proteins have been detected  
41  
42 independently by Rohmer *et al.* [29]. This overlap corresponds to a p-value of  $3.4 \cdot 10^{-6}$   
43  
44 under the assumption of a random choice of the 31 proteins.  
45

46 Finally, the estimated FDR of 5.7 % of the OPM-based analysis corresponds to 15  
47  
48 false positives out of the 257 proteins detected. Therefore even in the unlikely case  
49  
50 where all proteins common to Mascot are real positives, such that all 15 false positive  
51  
52 are clustered in the 31 additional proteins, at least 16 of these would be real new hits.  
53  
54  
55  
56  
57  
58  
59  
60

## 4 Discussion

### 4.1 Summary of results

In the present paper we describe an alignment approach that exploits both peptide mass and retention order for the identification of proteins. The main result of this work is that a combination of highly accurate MS, modern retention time prediction and robust scoring algorithm can be used to identify a large number of proteins, of the same order of magnitude as fragmentation-based methods. The scoring algorithm introduced here exploits both peptide mass and retention order for the identification of proteins. The new method tends to detect the same proteins as those detected in the fragmentation-based approach. It evaluates peptide peaks only within the context of the candidate proteins from which they arise.

Using this method we were able to identify 257 proteins from *Francisella tularensis*. This is more than twenty times the 12 proteins repeatedly detected by Palmblad *et al.* from membranes of *Yersinia pestis* [20] in the only work that we are aware of, that reports concrete identification results from fragmentation-free MS and predicted HPLC retention times. Relative to the total number of proteins in the samples (at most 1719 in *Francisella*, and 456 in the membrane fraction of *Yersinia pestis* [30]), we detect a much higher proportion of the proteins present in the sample (15% vs 2.6%) and are able to estimate a confidence level for these protein identifications.

## 4.2 Robustness of the scoring algorithm with regard to prediction errors

As our scoring algorithm only relies on prediction *order*, it is robust against differences between predicted and measured retention time. This relaxes the accuracy requirement on the retention time predictor used: Differences between actual and predicted retention times will not change the outcome as long as they do not alter the retention order.

Furthermore, if the predicted retention order for two peptides is inverted, both can still contribute to the OPM score: As long as they elute in the overlapping sections of their retention profiles, the peaks corresponding to these peptides can be found by the alignment, despite their wrong order.

Finally, a correct retention order used is only relevant when comparing peptides resulting from the digest of the same protein. An incorrect retention order for peptides from other proteins has no influence on the OPM score.

To test these ideas, we have simulated errors by adding Gaussian-distributed noise to each of the predicted peptide retention times and rerunning the analysis. We evaluated the effect of this modification on the estimated false discovery rate and on the number of proteins also found by fragmentation and Mascot. The magnitude of this noise is expressed in the same units as the hydrophobic index used by SSRCalc and described in Section **1.2.2**.

For a Gaussian noise with a standard deviation below 4 units, the FDR only slightly increased for all numbers of positive proteins. Similarly the ROC curve was only slightly modified. However when the standard deviation of the noise reached 5 units and more, the FDR increased to the levels obtained when the order of peptides was not taken into account and even beyond that level. The same degradation was seen on the ROC curve. For comparison purposes we computed the intervals between

1  
2  
3  
4  
5  
6  
7  
8 the predicted retention times of subsequent tryptic peptides belonging to the same  
9 proteins for all proteins in the database. Over the whole database, we found that the  
10 median of these intervals was 1.18 SSRCalc units and the 3rd quartile 2.84 units. This  
11 is consistent with the magnitude of the introduced noise for which the performance of  
12 our alignment-based detection method starts to be disrupted.  
13  
14  
15  
16

### 17 18 **4.3 Impact of the ion pairing reagent**

19  
20  
21  
22 Commonly used ion pairing agents for the HPLC of peptides include Trifluoroacetic  
23 Acid (TFA) and Formic Acid (FA). The retention time predictor used here [21, 25]  
24 had been optimized by its authors for TFA on 100Å pore-size columns. On the other  
25 hand, the data that we used have been obtained with FA. Although the authors  
26 emphasized that the performance of their predictor was lower in FA than in TFA (the  
27 correlation coefficient between predictions and measures was only  $R^2 = 0.94$  with FA  
28 as opposed to a typical  $R^2 = 0.99$  with TFA on different columns), we still obtained an  
29 improvement of protein identification with the retention order in TFA. We interpret  
30 this result as evidence for the robustness of the OPM method.  
31  
32  
33  
34  
35  
36  
37  
38  
39

### 40 41 **4.4 Impact of sample complexity and database**

#### 42 43 44 **size**

45  
46  
47  
48 The performance of most protein identification approaches depends on the sample  
49 complexity and the size of the database: The more complex the sample, and the larger  
50 the database used for identification, the less one will be able to infer from MS peaks.  
51 We believe this to be true for the OPM approach as well. The examples shown here  
52 were obtained with samples of medium complexity, namely soluble proteins of whole  
53 cell lysate, from a prokaryote with a relatively small proteome: With its estimated  
54 1719 proteins, *Francisella tularensis* has a simple proteome: 262 (13.9%) of the 1886  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 entries listed in the Integr8 [1] database of complete genomes (version 110, May 26<sup>th</sup>,  
9 2010) include a number of proteins equal or smaller to that value. Similarly, 50% of  
10 the genomes in Integr8, V 110 have less than 3187 proteins.  
11  
12

13  
14 It is noteworthy that other strategies that require limited fragmentation data were  
15 also developed on target proteomes of relatively low complexity. For instance, the  
16 AMT strategy was developed on *Deinococcus radiodurans*, which has a proteome of  
17 3085 proteins [11] and the *Yersinia pestis* membrane fraction used by Palmblad [20]  
18 contains only 456 proteins.  
19  
20  
21  
22

## 23 24 **4.5 Post-translational modifications and missed** 25 **clivages.** 26 27 28

29  
30 Individual amino acids in proteins are often modified from their canonical form after  
31 protein synthesis. These post-translational modifications (PTMs) imply modifications  
32 in the measured masses and retention times. They can either affect every instance of  
33 the relevant amino acid (fixed PTM) or only a subset of *a priori* unknown size  
34 (variable PTM). Similarly, some predicted enzymatic cleavages may not occur  
35 experimentally (missed cleavage).  
36  
37  
38  
39  
40  
41  
42

43 As presented here, the OPM approach does not accommodate PTMs or missed  
44 cleavages, but it can easily be extended for both. Conceptually, fixed missed  
45 cleavages merely lead to a different set of peptides whose retention times can also be  
46 predicted. Variable missed cleavages are more problematic: because of the global  
47 nature of the scoring scheme, the match of several peptides could be affected by the  
48 modification of only one of them and the optimal alignment needs to be recomputed.  
49  
50  
51  
52

53  
54 The introduction of a PTM induces changes in both peptide mass and retention  
55 time. Any fixed PTM taken into account by the retention time predictor (as is the case  
56 here for the alkylation by iodoacetamide) only translates into a different set of ordered  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 peptides as input for the OPM method. In the case of variable PTMs, two possibilities  
9  
10 (with and without the PTM) must be considered individually. As in other methods, the  
11  
12 number of cases to consider will increase exponentially with the number of  
13  
14 independent, variable PTMs in a given protein.

15  
16 The case of phosphorylations is of particular importance. According to Steen *et al.*  
17  
18 [31], the retention time of peptides is, in many cases, not drastically modified upon  
19  
20 phosphorylation, particularly if the number of basic residues in the peptide and the  
21  
22 positive charges that they carry can counterbalance the additional negative charge  
23  
24 induced by the phosphorylation. This would happen with monophosphorylated tryptic  
25  
26 peptides which almost always have a basic residue (lysine or arginine) at their C-  
27  
28 terminal end. Therefore as the OPM method is robust to minor errors in the retention  
29  
30 time it could readily accommodate phosphorylation using existing retention time  
31  
32 prediction methods.

## 33 34 35 36 **4.6 Incorporating additional information** 37 38

39  
40 The interest of a protein identification method based on HPLC and single MS could  
41  
42 be questioned in view of the increasing availability of mass spectrometers with  
43  
44 fragmentation and MS/MS capacity.

45  
46 However, even with the current generation of fast instruments, it is still not  
47  
48 possible to fragment all the peptide ions from the first MS. Furthermore, not all the  
49  
50 precursor ions selected for fragmentation provide spectra leading to the identification  
51  
52 of peptide sequences. Therefore, the primary scans constitute a more complete  
53  
54 representation of the peptides present in the samples than the series of assigned  
55  
56 fragmentation spectra. The present method allows the use of data from this faithful  
57  
58 representation of the digest peptides.

59  
60 Finally, the OPM approach can be used to incorporate additional information into



1  
2  
3  
4  
5  
6  
7  
8 fragment-based analysis: For instance, in HPLC-MS/MS experiments, some proteins  
9  
10 only give rise to a single fragmentation spectrum, which is frequently not considered  
11  
12 as sufficient for unambiguous identification. Additional evidence about the presence  
13  
14 of such a protein could be obtained by searching for peaks from its other peptides in  
15  
16 the primary spectra, even if these peptides have not been fragmented. The scoring  
17  
18 method described here can be modified for this purpose: If a peptide has been  
19  
20 detected by fragmentation analysis, the alignment score for matching with this  
21  
22 peptide is increased strongly above the default matching score. The rest of the  
23  
24 alignment is executed as described above, rewarding matches of additional predicted  
25  
26 peptides with the primary spectra. With the appropriate score settings, the resulting  
27  
28 optimal alignment will incorporate the peptide recognized after fragmentation and the  
29  
30 protein identification will thus be based on evidence from both primary and  
31  
32 secondary spectra. Alternatively mass and retention time information from one  
33  
34 experiment could be used in further experiments to guide the choice of precursor  
35  
36 peptides, and thus provide more informative fragmentation steps.

## 37 38 39 **5 Conclusions**

40  
41  
42 Our main conclusion is that highly accurate mass measurements, modern retention  
43  
44 time prediction, and a new scoring algorithm can be combined into a fragmentation-  
45  
46 free approach that is able to identify hundreds of proteins in a medium-complexity  
47  
48 sample. This demonstrates for the first time that, supported by suitable informatics  
49  
50 approaches, fragmentation-free mass spectrometry may, in the future and depending  
51  
52 on sample complexity, deliver protein identification performance on a level similar to  
53  
54 that of fragmentation-based approaches. In this respect, new instruments enable this  
55  
56 possibility, for example the Thermo Exactive mass spectrometer  
57  
58 (<http://www.thermo.com/Exactive>).

59  
60 The OPM method introduced here provides a novel approach of using the

1  
2  
3  
4  
5  
6  
7  
8 information on peptide retention time for the identification of peptides by HPLC-MS.  
9  
10 Contrary to existing methods [11] it does not require previous retention time  
11 measurement. Furthermore, it relies on the retention *order* of peptides, and not on  
12 their exact retention times. Therefore, it can use predictors not trained with data  
13 obtained in exactly the same conditions and on the same HPLC system used for the  
14 analysis. In the practical application which we studied here, the OPM approach  
15 identified a significant fraction of the proteins found by a fragmentation-based  
16 approach, and some additional proteins.  
17  
18  
19  
20  
21  
22

23  
24 Finally, the alignment method underlying the OPM approach may also be used in  
25 combination with fragmentation data in at least two ways: 1) The alignment for one  
26 specific protein can include one peptide recognized from MS/MS spectra and other  
27 peptides from the same protein matching peaks in primary scans. 2) The alignment  
28 results (mass and retention time of the peptides) can direct the choice of ions which  
29 will be submitted to fragmentation in further LC/MS-MS experiments with the same  
30 samples.  
31  
32  
33  
34  
35  
36  
37

## 38 **6 Acknowledgments**

39  
40  
41  
42 Work of P. Clote was partially funded by the Digiteo Foundation, as well as by the  
43 National Science Foundation grants DBI-0543506 and DMS-0817971. Any opinions,  
44 findings, and conclusions or recommendations expressed in this material are those of  
45 the authors and do not necessarily reflect the views of the National Science  
46 Foundation. Frank Rügheimer was supported by a grant of the European Union (FP6,  
47 BaSysBio, grant LSHG-CT-2006-037469). The authors would like to thank Sébastien  
48 Li-Thiao-Té for discussions and suggestions and Thomas Rolland for comments.  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 7 Conflict of interest

The authors declare no conflict of interest.

## References

- [1] Kersey, P., Bower, L., Morris, L., Horne, A., *et al.* Integr8 and genome reviews: integrated views of complete genomes and proteomes. *Nucleic Acids Res* 2005. 33, D297–D302.
- [2] Mann, M., Hojrup, P., Roepstorff, P. Use of mass spectrometric molecular weight information to identify proteins in sequence databases. *Biol Mass Spectrom* 1993. 22, 338–345.
- [3] Pappin, D. J., Hojrup, P., Bleasby, A. J. Rapid identification of proteins by peptide-mass fingerprinting. *Curr Biol* 1993. 3, 327–332.
- [4] Perkins, D. N., Pappin, D. J., Creasy, D. M., Cottrell, J. S. Probability-based protein identification by searching sequence databases using mass spectrometry data. *Electrophoresis* 1999. 20, 3551–3567.
- [5] Zhang, W., Chait, B. T. Profound: an expert system for protein identification using mass spectrometric peptide mapping information. *Anal Chem* 2000. 72, 2482–2489.
- [6] Eng, J. K., McCormack, A. L., III, J. R. Y. An approach to correlate tandem mass spectral data of peptides with amino acid sequences in a protein database. *Journal of the American Society for Mass Spectrometry* 1994. 5, 976 – 989.
- [7] McCormack, A. L., Schieltz, D. M., Goode, B., Yang, S., *et al.* Direct analysis and identification of proteins in mixtures by lc/ms/ms and database searching at the low-femtomole level. *Anal Chem* 1997. 69, 767–776.
- [8] Baczek, T., Kaliszan, R. Predictions of peptides' retention times in reversed-

- 1  
2  
3  
4  
5  
6  
7  
8 phase liquid chromatography as a new supportive tool to improve protein  
9 identification in proteomics. *Proteomics* 2009. 9, 835–847.
- 10  
11 [9] Norbeck, A. D., Monroe, M. E., Adkins, J. N., Anderson, K. K., *et al.* The utility  
12 of accurate mass and lc elution time information in the analysis of complex  
13 proteomes. *J Am Soc Mass Spectrom* 2005. 16, 1239–1249.
- 14  
15  
16  
17 [10] Meek, J. L. Prediction of peptide retention times in high-pressure liquid  
18 chromatography on the basis of amino acid composition. *Proc Natl Acad Sci U*  
19 *SA* 1980. 77, 1632–1636.
- 20  
21  
22  
23 [11] Strittmatter, E. F., Ferguson, P. L., Tang, K., Smith, R. D. Proteome analyses  
24 using accurate mass and elution time peptide tags with capillary lc time-of-flight  
25 mass spectrometry. *J Am Soc Mass Spectrom* 2003. 14, 980–991.
- 26  
27  
28  
29 [12] Guo, D., Mant, C. T., Taneja, A. K., Rodges, R. S. Prediction of peptide  
30 retention times in reversed-phase high-performance liquid chromatography ii.  
31 *Journal of Chromatography A* 1986. 359, 519 – 532.
- 32  
33  
34  
35 [13] Palmblad, M., Ramström, M., Markides, K. E., Håkansson, P., Bergquist, J.  
36 Prediction of chromatographic retention and protein identification in liquid  
37 chromatography/mass spectrometry. *Anal Chem* 2002. 74, 5826–5830.
- 38  
39  
40  
41 [14] Mant, C. T., Zhou, N. E., Hodges, R. S. Correlation of protein retention times in  
42 reversed-phase chromatography with polypeptide chain length and  
43 hydrophobicity. *J Chromatogr* 1989. 476, 363–375.
- 44  
45  
46  
47 [15] Petritis, K., Kangas, L. J., Ferguson, P. L., Anderson, G. A., *et al.* Use of  
48 artificial neural networks for the accurate prediction of peptide liquid  
49 chromatography elution times in proteome analyses. *Anal Chem* 2003. 75,  
50 1039–1048.
- 51  
52  
53  
54 [16] Petritis, K., Kangas, L. J., Yan, B., Monroe, M. E., *et al.* Improved peptide  
55 elution time prediction for reversed-phase liquid chromatography-ms by  
56 incorporating peptide sequence information. *Anal Chem* 2006. 78, 5026–5039.
- 57  
58  
59  
60

- 1  
2  
3  
4  
5  
6  
7  
8 [17] Pfeifer, N., Leinenbach, A., Huber, C. G., Kohlbacher, O. Statistical learning of  
9 peptide retention behavior in chromatographic separations: a new kernel-based  
10 approach for computational proteomics. *BMC Bioinformatics* 2007. 8, 468.  
11  
12 [18] Pfeifer, N., Leinenbach, A., Huber, C. G., Kohlbacher, O. Improving peptide  
13 identification in proteome analysis by a two-dimensional retention time filtering  
14 approach. *J Proteome Res* 2009. 8, 4109–4115.  
15  
16 [19] Krokhin, O. V. Sequence-specific retention calculator. algorithm for peptide  
17 retention prediction in ion-pair rp-hplc: application to 300- and 100-Å pore size  
18 c18 sorbents. *Anal Chem* 2006. 78, 7785–7795.  
19  
20 [20] Palmblad, M., Ramström, M., Bailey, C. G., McCutchen-Maloney, S. L., *et al.*  
21 Protein identification by liquid chromatography-mass spectrometry using  
22 retention time prediction. *J Chromatogr B Analyt Technol Biomed Life Sci* 2004.  
23 803, 131–135.  
24  
25 [21] Krokhin, O. V., Ying, S., Cortens, J. P., Ghosh, D., *et al.* Use of peptide  
26 retention time prediction for protein identification by off-line reversed-phase  
27 hplc-maldi ms/ms. *Anal Chem* 2006. 78, 6265–6269.  
28  
29 [22] Oyston, P. C. F., Sjostedt, A., Titball, R. W. Tularemia: bioterrorism defence  
30 renews interest in francisella tularensis. *Nat Rev Microbiol* 2004. 2, 967–978.  
31  
32 [23] Guina, T., Radulovic, D., Bahrami, A. J., Bolton, D. L., *et al.* MglA regulates  
33 francisella tularensis subsp. novicida (francisella novicida) response to  
34 starvation and oxidative stress. *J Bacteriol* 2007. 189, 6580–6586.  
35  
36 [24] Matthiesen, R. Extracting monoisotopic single charge peaks from liquid  
37 chromatography electrospray ionization mass spectrometry. *Methods Mol Biol*  
38 2007. 367, 37–48.  
39  
40 [25] Spicer, V., Yamchuk, A., Cortens, J., Sousa, S., *et al.* Sequence-specific  
41 retention calculator. a family of peptide retention time prediction algorithms in  
42 reversed-phase hplc: applicability to various chromatographic conditions and  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 columns. *Anal Chem* 2007. 79, 8762–8768.

- 9  
10 [26] Koenker, R. *Quantile Regression*. Econometric Society Monographs.  
11 Cambridge University Press, New York, NY, 2005.  
12  
13 [27] Koenker, R. *quantreg: Quantile Regression*, 2008. R package version 4.24.  
14  
15 [28] Elias, J. E., Gygi, S. P. Target-decoy search strategy for increased confidence in  
16 large-scale protein identifications by mass spectrometry. *Nat Methods* 2007. 4,  
17 207–214.  
18  
19 [29] Rohmer, L., Guina, T., Chen, J., Gallis, B., *et al.* Determination and comparison  
20 of the francisella tularensis subsp.novicida u112 proteome to other bacterial  
21 proteomes. *J Proteome Res* 2008. 7, 2016–2024.  
22  
23 [30] Pieper, R., Huang, S.-T., Clark, D. J., Robinson, J. M., *et al.* Integral and  
24 peripheral association of proteins and protein complexes with yersinia pestis  
25 inner and outer membranes. *Proteome Sci* 2009. 7, 5.  
26  
27 [31] Steen, H., Jebanathirajah, J. A., Rush, J., Morrice, N., Kirschner, M. W.  
28 Phosphorylation analysis by mass spectrometry: myths, facts, and the  
29 consequences for qualitative and quantitative measurements. *Mol Cell*  
30 *Proteomics* 2006. 5, 172–181.  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 8 Figures and table legends

Figure 1: Principle of the identification method: for each protein the predicted tryptic peptides are sorted according to their predicted retention time (bottom). They are then matched to ions from the spectra in the same order. Each trapezoid represents one experimental spectrum, and the masses of the peptides are symbolized by the height of the lines. Intensities of the peaks are not considered. Solid lines: Matching peptides; Dashed lines: Non matching peptides.

Figure 2: This dynamic programming algorithm iteratively computes the score  $s_{i,j}$  of the best alignments between a prefix  $(p_1, \dots, p_j)$  of the list of peptide masses and a prefix  $(S_1, \dots, S_i)$  of the list of spectra. Here the values of the *no-peptide score* and the *no-spectrum score* are 0 and  $-1$  respectively.

Figure 3: Distribution of number of matched peptides for proteins from *Francisella tularensis subsp. novicida* with one HPLC-MS experiment. Each point represents a protein from the sequence database. For visualization only, the positions of the points are shifted by a small gaussian random amount ( $\sigma=0.5$ ). The lines represent the best linear fits for the data (see text). Abscissa : number of peptides longer than 4 amino acids per protein; Ordinate: number of matched peptides.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Figure 4: Effect of incorporating the order of the peptides. A : Values of the estimated false discovery rate (FDR) for different numbers of proteins identified above the threshold, with or without taking the peptide order into account. B : Receiver Operating Characteristic (ROC) curves for the comparison of the novel method with results from peptide fragmentation and Mascot. The ROC curves are obtained using the output from Mascot as the reference set.

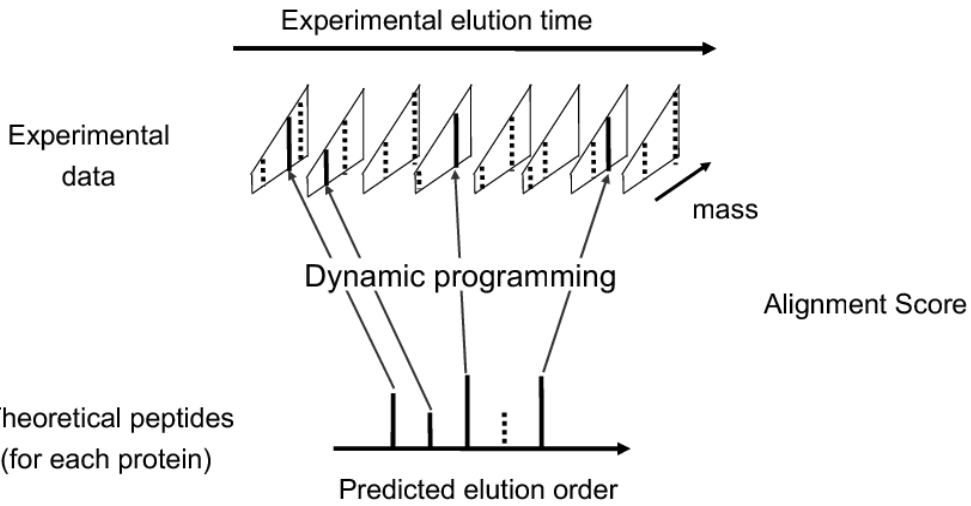
Table 1: OPM scores and quantile assignment of *F. tularensis* proteins. For each protein this table lists the number of tryptic peptides longer than 4 amino acids and the number of these peptides matched to the experimental MS1 peaks by the alignment procedure. In the last column the quantile assignment of the protein is shown. This corresponds to the position of the alignment score in the distribution of all the scores for decoy proteins with the same number of tryptic peptides.



1  
2  
3  
4  
5  
6  
7  
8 Table 2: Number and proportion of total proteins detected for several quantile  
9  
10 limits. Top: with experimental peptides matching in the predicted order; Bottom:  
11  
12 without ordering the peptides. Real proteins: number of real proteins detected (portion  
13  
14 of total proteins); Random proteins: number of randomized proteins from the decoy  
15  
16 database detected; FDR estimate: estimated proportion of absent proteins from the  
17  
18 protein with scores above the quantile limit; Frag. overlap: detected proteins also  
19  
20 identified from the fragmentation data with the Mascot software. The total number of  
21  
22 proteins in the database is 1719 and 436 of them were identified by Mascot from the  
23  
24 fragmentation spectra.  
25  
26  
27  
28  
29  
30

31 Table 3: Proteins identified by fragmentation and Mascot. This table lists the  
32  
33 accession numbers and Mascot scores for the 436 proteins detected by Mascot on the  
34  
35 basis of the fragmentation (MS-MS)spectra. Based on the comparison with a decoy  
36  
37 data base this number of proteins was retained in order to obtain a false discovery rate  
38  
39 of 5.7%.  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



15  
16  
17  
18 *Algorithm OrderedPeptideMatchScore*  
19

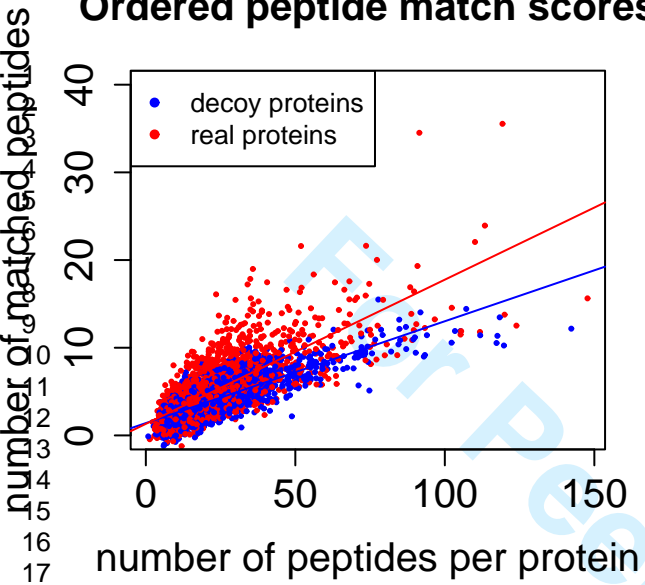
20 Input:

- A protein  $P$
- A list  $(S_1, \dots, S_m)$  of decharged, deisotoped spectra
- A list of masses  $(p_1, \dots, p_n)$  of the tryptic peptides of  $P$ , ordered by predicted elution time
- A *peak match threshold*  $\Delta$
- Constant *no-peptide score*  $nps$ , *no-spectrum score*  $nss$ , and a *peptide-spectrum match score*  $ps(q)$

- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29 1. For  $i=1, \dots, m$
- 30 2. For  $j=1, \dots, n$
- 31 3. If  $S_i$  has a peak in the interval  $(p_j - \Delta, p_j + \Delta)$ , set  $ps = +1$ ; otherwise set  
32  $ps = -1$ .
- 33 4. Compute,  $s_{i,j} = \max(s_{i-1,j-1} + ps, s_{i-1,j} + nps, s_{i,j-1} + nss)$  under the convention that  
34  $s_{i,j} = 0$  if  $i=0$  or  $j=0$

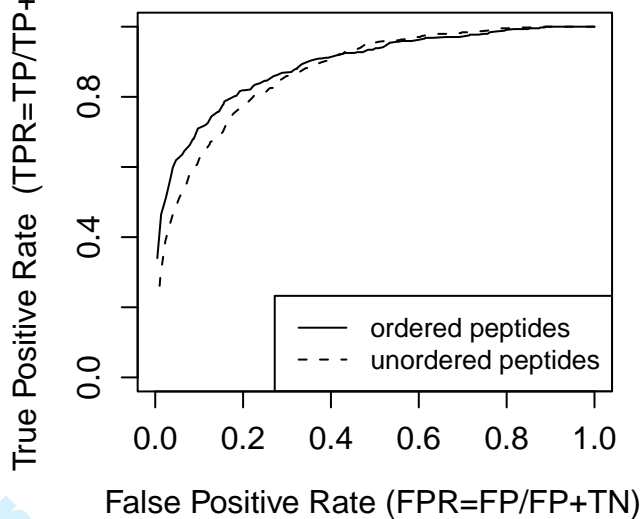
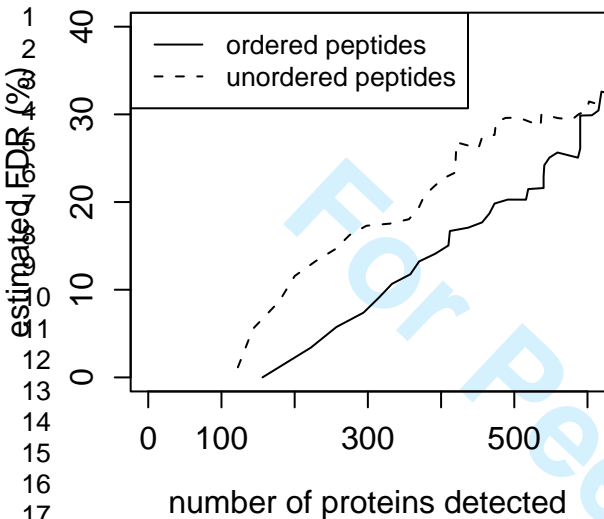
35 Output:  $s_{m,n}$   
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

# Ordered peptide match scores



**A: False Discovery Rate**

**B: ROC curves, detector vs Mascot**



Ordered peptide match score

Quantile	0.80	0.90	0.95	0.99
Real proteins	721 (41.9%)	531 (30.9%)	407 (23.7%)	257 (14.9%)
Random proteins	372 (21.6%)	174 (10.1%)	93 (5.4%)	15(0.91%)
FDR estimate	34.60%	24.80%	16.90%	5.7%
Frag. overlap	376	337	295	226

Unordered peptide match score

Quantile	0.80	0.90	0.95	0.99
Real proteins	701 (40.8%)	479 (27.9%)	357 (20.8%)	181 (10.5%)
Random proteins	350 (20.4%)	168 (9.8%)	85 (4.9%)	16 (0.9%)
FDR estimate	36.3%	28.8%	20.0%	8.6%
Frag. overlap	367	301	256	158