

# BIOSPEAN: A Freeware Tool for Processing Spectra from MALDI Intact Cell/Spore Mass Spectrometry

Martin Raus and Marek Šebela\*

Department of Protein Biochemistry and Proteomics, Centre of the Region Haná for Biotechnological and Agricultural Research, Faculty of Science, Palacký University, Šlechtitelů 11, CZ-783 71 Olomouc, Czech Republic

## Abstract

Here we introduce the software BIOSPEAN, which is applicable for processing of peptide/protein profiles obtained from MALDI intact cell/spore mass spectrometry. It has been developed as a web application using common technology. The BIOSPEAN automatically finds peaks in a mass spectrum. The peak detection principle involves a local scanning of intensity values around individual  $m/z$  positions. To cope with the level of noise, a threshold signal-to-noise ratio is adjusted for filtering relevant results. Based on the detected peak pattern, a similarity search in a custom spectral database can subsequently be performed for sample identification. When a spectrum is analyzed by comparison with a database entry, a percentage score value is calculated from the number of identical peak positions found (they are assigned with an adjustable mass tolerance) divided by the number of all detected peaks in the analyzed spectrum. Each two spectra may also be compared in the opposite way and both score values averaged. The database search results are finally sorted in a table.

**Keywords:** Comparison; Database search; Mass spectrometry; Microorganisms; Peak; Software

**Abbreviations:** ESI: Electrospray Ionization; FA: Ferulic Acid; IC/IS: Intact Cell or Intact Spore; MALDI-TOF: Matrix-Assisted Laser Desorption/Ionization Time-of-Flight; MS: Mass Spectrometry; MSP: Main Spectrum; SA: Sinapinic Acid

## Introduction

Since the discovery of soft ionization techniques at the end of the 1980s, mass spectrometry (MS) has become an important tool in biological research [1]. Mass spectrometers with electrospray ionization (ESI) or matrix-assisted laser desorption/ionization (MALDI) are indispensable in disciplines such as proteomics, metabolomics or microbial identification [2,3]. The instruments are provided by their vendors with a licensed software, which is utilized by users for the acquisition (including instrument calibration) and saving of mass spectra, processing of the acquired spectra (smoothing, baseline correction, peak picking, peaklist building etc.), spectra comparison and searches against spectral or sequence databases. Moreover, there is a freeware or open-source software available for that purpose such as the mMass [4]. The peak pattern represents a characteristic feature of each mass spectrum. Thus its determination is challenging for any processing software. A good peak-picking algorithm must remove noise without removing weak signals [5]. Moreover, the identification of peaks has to be done with a high precision to recognize overlapping peaks and filter true peaks from a mixed set of true and false peaks. The request for archiving of spectra in the form of a library or database and their sorting appears immediately when there is a need to work with more than a single spectrum. Mass spectral databases are essential namely for the rapid characterization and identification of microorganisms [6]. Assigning of an experimental spectrum to its counterpart in the database is done through the highest level of matching of specific peaks constituting a fingerprint.

Two comprehensive software solutions for biotyping of microorganisms were introduced in the 2000s by AnagnosTec (Zossen, Germany) and Bruker Daltonik (Bremen, Germany). AnagnosTec developed a software with database [7], which has been named SARAMIS (Spectral ARchiving And Microbial Identification

System) and commercialized. In 2010, the database was acquired by the company bioMérieux (Marcy l'Etoile, France) and currently it is provided as a part of the VITEK MS system together with a MALDI-TOF (the abbreviation TOF stands for "time-of-flight") instrument supplied by the Shimadzu Corporation [8]. The SARAMIS database stores ten thousands of single fingerprint spectra of different microbial isolates representing more than 2,000 species and 500 genera [9]. The identification of microorganisms is done through the concept of so-called SuperSpectra. This name refers to MS data that are acquired with well-characterized microorganisms and highlight biomarker features at the level of genus, species and strain. Each "super" spectrum is generated from isolates of one species originating from different hospitals, reference laboratories or culture collections. The identification is based on a percentage of confidence returned from mutual comparisons with reference spectra. Empirical threshold "confidence values" of 90 and 70 % are recommended for species- or genus-level identification, respectively [10]. Besides this function, the SARAMIS platform also allows construction of dendrograms to show taxonomic relationships among strains.

The biotyping system by Bruker consists of a MALDI-TOF mass spectrometer and the MALDI Biotyper software. Also in this case, the identification is done by comparison of the respective sample spectrum with a library of reference spectra. The pattern recognition and matching algorithm used considers positions, intensities and

\*Corresponding author: Marek Šebela, Department of Protein Biochemistry and Proteomics, Centre of the Region Haná for Biotechnological and Agricultural Research, Faculty of Science, Palacký University, Šlechtitelů 11, CZ-783 71 Olomouc, Czech Republic, Tel: +420 585634927; Fax: +420 585634933; E-mail: [marek.sebela@upol.cz](mailto:marek.sebela@upol.cz)

Received October 07, 2013; Accepted November 25, 2013; Published November 28, 2013

Citation: Raus M, Šebela M (2013) BIOSPEAN: A Freeware Tool for Processing Spectra from MALDI Intact Cell/Spore Mass Spectrometry. J Proteomics Bioinform 6: 282-287. doi:10.4172/jpb.1000292

Copyright: © 2013 Raus M, et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

frequencies of peaks. The reference spectra are termed “Main Spectra” (abbreviated as MSPs) and each of them comes from a multiple analysis of a single and defined strain [11]. In 2013, the MALDI Biotyper system covered a broad range of more than 4,600 microbial isolates from gram-negative bacteria, gram-positive bacteria, yeasts, multicellular fungi and mycobacteria (more than 2,000 species). Database search results are sorted according to score values [12]. The range from 2.300 to 3.000 (green label) represents a highly probable species identification, score values from 2.000 to 2.299 (still green label) secure genus identification when assignment to a species becomes probable. No reliable identification is achieved with score values ranging from 0.000 to 1.699 (red label) whereas values appearing between 1.700 and 1.999 suggest only a probable species identification (yellow label). In addition to data processing and identification via database search, the software allows obtaining further information such as the construction of taxonomic trees, principal component analysis etc.

In addition to Bruker and bioMérieux, the Andromas software and database for microbial identification using MALDI-TOF MS [13] has been offered since 2012 by the homonymous company located in Paris, France. However, this product (currently integrated with a MALDI-TOF mass spectrometer) seems to be confined to the French market only. The Andromas strategy of identification is based on registration of a limited number of species-specific peaks in the acquired mass spectrum (suggested arbitrarily as peaks with a relative intensity higher than 7% when compared to the base peak). Database search results with sample spectra are displayed in percentage values representing similarity to reference spectra. The term “good identification” for example refers to a similarity value of at least 65% with a difference between the first two assigned species of at least 10% [13]. To complete this overview, the software MicrobeLynx by Waters Corporation (Manchester, UK) should be mentioned, which appeared at the beginning of commercialization of MALDI-based microbial identification systems [14]. That time the software utilized only a limited library of reference spectra (circa 3,500 entries) and has not achieved a wide distribution ever since.

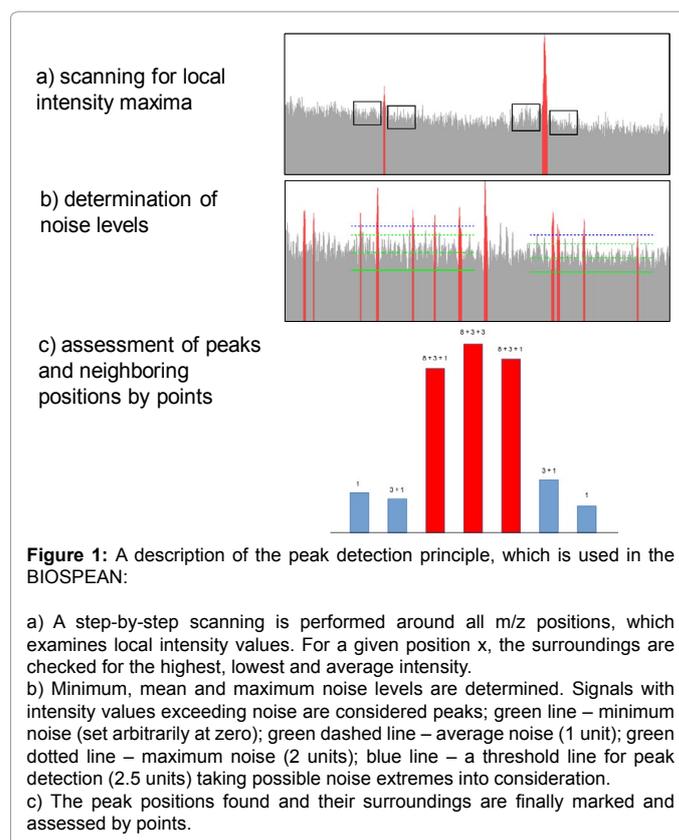
The commercial software is distributed with a single user license or the customers may purchase a multi-user license, which is of course reflected in contract conditions and the final price of the whole system. In our laboratory, we are involved in biochemical and proteomic applications of MS including analyses of fungi using intact cell/spore mass spectrometry (IC/IS MS) performed on MALDI-TOF instruments [15]. The initial idea behind the software BIOSPEAN, which is introduced here, was to develop a freeware tool (preferably online accessible) for processing fingerprint spectra from IC/IS MS i.e. to perform peak picking, building peak lists, construction of a spectral database and data sharing. The most important benefit should reside in the possibility of doing sample spectra comparisons with database entries for species identification. Furthermore, options such as adding user-defined functions and changing parameters were highly desirable. The software now runs in a test mode. Recently we showed its applicability for the analysis of isomeric *O*-glycopeptides, which was based on construction of a library of virtual tandem mass spectra covering theoretical fragmentation patterns with an indicative distribution of reporter  $\gamma$ -ions [16].

## Features of the BIOSPEAN

The BIOSPEAN (see at <http://biochemie.upol.cz/index.php/cs/vyzkum/odkazy>, under the heading Software) is a web application,

which allows performing tasks connected with a processing, storage and comparison of mass spectra. It has been developed primarily for the use in microbial identification using IC/IS MS. Prior to analysis, each mass spectrum is uploaded as a data file (such as a TXT- or CSV-formatted file) obtained via the respective acquisition software on a mass spectrometer. For that purpose, the BIOSPEAN relies on a common open source database system [17]. After uploading a new spectrum, its analysis and preprocessing is done. This includes a detection of peaks based on a local scanning over all  $m/z$  values ( $x$ -coordinates) for extreme intensity values ( $y$ -coordinates). The peak positions found are then assessed by points that reflect peak widths (Figure 1). For each detected peak, the respective position is assigned 8 points. The neighboring positions at a distance of 1 unit are assigned 3 points and those at a distance of 2 units receive 1 point each. If there are two and more neighboring peaks (or broad peaks), the assigned points are summed up (Figure 1). This is important for the subsequent comparison of spectra as it eliminates small mass shifts to some extent (up to 2 mass units).

The user may further modify each preprocessed spectrum. Besides adding information about the spectrum itself (a title, description, remarks, categorization), it is also possible to change the precision of peak detection. Mass spectra differ in their quality represented by the noise level and variability in peak intensities. For that reason, there is no universal setting of sensitivity parameters applicable for both peak recognition and noise removal. An average value of sensitivity is used at the beginning, which is rather empiric and comes from a repeated optimization with experimental data (Figure 1). However, this value does not need to be optimal for all analyzed spectra and thus the user is provided by the possibility of an individual adjustment of the sensitivity



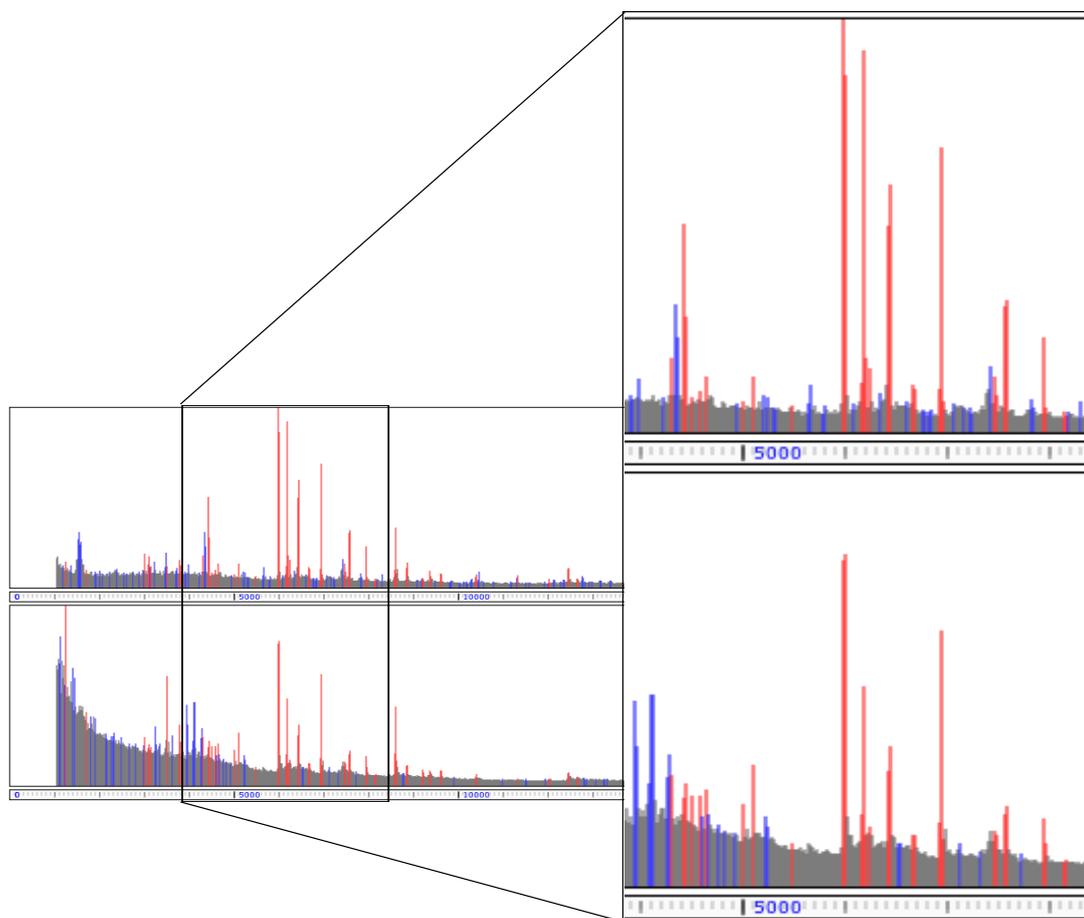
parameter for each spectrum in order to obtain the most optimal output. For that purpose, the user is guided by statistical graphs, which show information about the distribution of noise and number of peaks resulting from the adjusted detection parameters.

When this procedure still is not helpful (it becomes very problematic namely in the case of spectra with an irregular distribution of noise or background), peaks may be labeled manually. Either those local extreme intensity values are additionally selected by the user that are considered real peaks (certainly if they were not recognized automatically) or an unlabeled can be done of false peaks. This kind of a three-step detection (where the first step is automatic, it is followed by the optional sensitivity adjustment for peak recognition and finally completed by the manual correction) allows to determine all peak positions reliably.

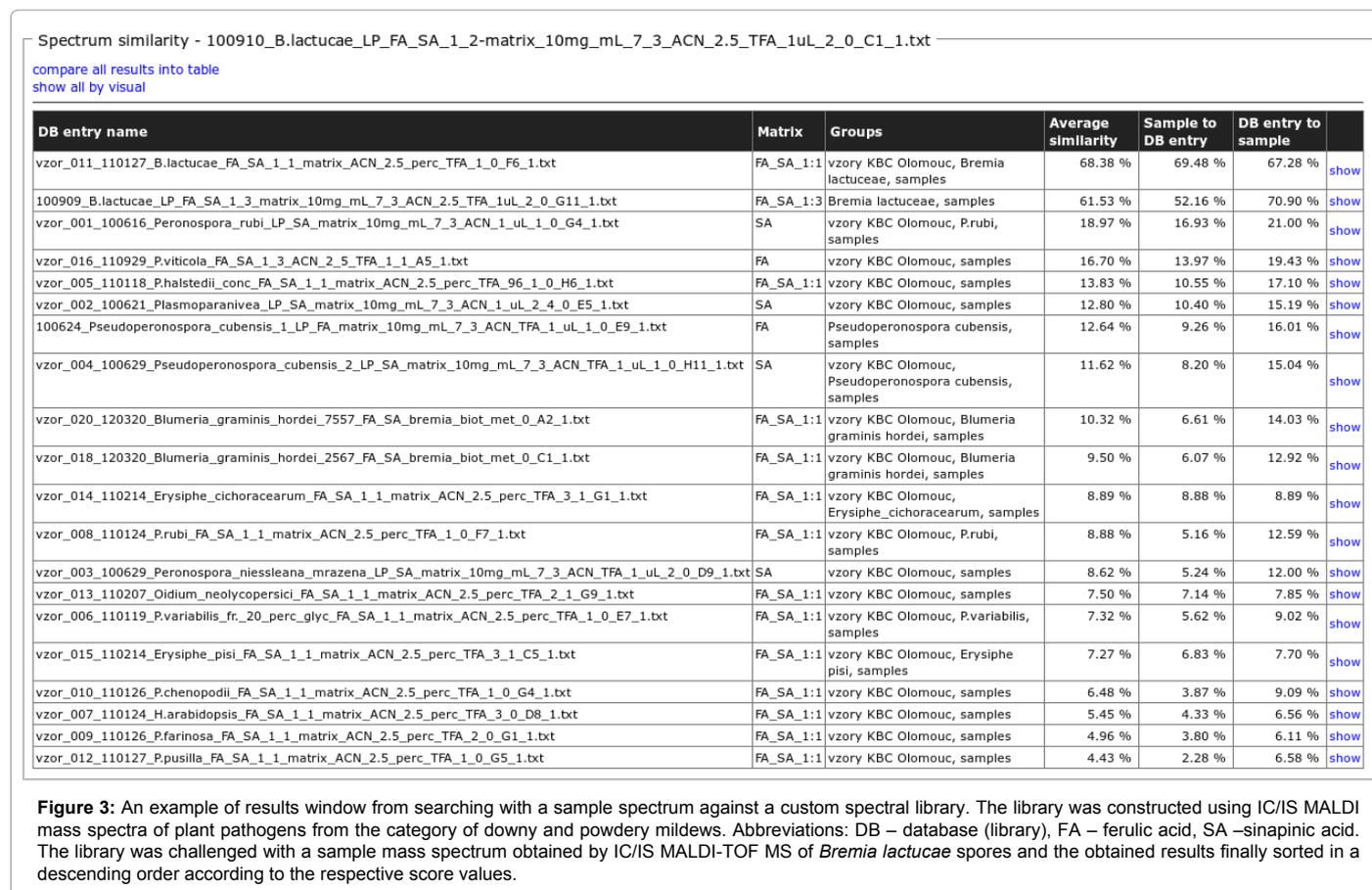
Since the spectrum has been placed in the database and assessed by points, it becomes available for a comparison with another spectrum (e.g. with a reference database entry). Any two selected spectra are compared on the basis of matching peak positions (Figure 2). A score value is calculated, which indicates how many percents of peaks in the analyzed spectrum possess their counterparts in the reference spectrum. Such a comparison is obviously asymmetric. When a spectrum A is compared with a spectrum B, the score value differs from that coming from a reversed comparison. Thus there are always two

values displayed as well as their average for every pair of spectra (A against B and B against A). Sorting of results for data interpretation then depends on user's decision.

The application BIOSPEAN is able to calculate a number of similarity values for a single analyzed spectrum resulting from its comparisons with all spectra in a user-made database. The results are listed in a table (Figure 3). The list of values can be sorted according to user's request e.g. using the average score value, according to the highest or lowest score, based on the best ratio of A/B and B/A or considering a complex-view coefficient, which takes into account both the ratio of similarity values and their average (in this case a ratio of 80/20 is less valuable than a ratio of 55/45). Figure 3 refers to a small library of reference IC/IS MALDI mass spectra representing various phytopathogenic fungal microorganisms studied in our laboratory. The library was challenged with a sample mass spectrum obtained by IC/IS MALDI-TOF MS of *Bremia lactucae* spores (see ref. [15] for further information). The identification was unambiguous. It is worth mentioning that the analyzed spectrum was measured under a different experimental setup (matrix composition) than the two *Bremia* spectra included in the reference library. For that reason, a score value of 70% was obtained whereas more convincing values above 90% are possible to reach in a standardized procedure. When the commercial software Biotyper 2.0 was used for a performance comparison, the library was first converted to a group of MSPs ("Main Spectra"). Then the analyzed



**Figure 2:** An example of spectra comparison window displayed in the software BIOSPEAN. The two compared mass spectra were acquired by IC/IS MS of fungal phytopathogen samples on a MALDI-TOF mass spectrometer. Identical peak positions found by the BIOSPEAN processing are displayed in red.



mass spectrum was subjected to a search in the group of MSPs. The obtained results are shown in Supplementary Figure. When compared with Figure 3, the same identification at the species level is obvious. Both lists of assigned library entries sorted in a descending order according to the respective score values are not completely identical (this can be explained by the different approaches for score calculations that are used in each individual software) but they are highly similar.

The compared spectra are visualized to highlight matched peaks as well as differences. A detailed view is available only for a simple comparison of two spectra (Figure 2). For a multiple comparison, less detailed view is offered with a vertical arrangement of panels. In both cases, matched peaks are distinguished from the others by color coding. When performing the multiple comparison, those peaks can also be visualized, which appear in a majority of experimental spectra but are missing in the others because of noise or an error originating from a wrong sample preparation or data acquisition on the instrument.

This is important for constructing virtual spectra used for example as species-specific or strain-specific references for bacteria or fungi. One may for example visualize all peaks that are present in 80 % of spectra from a selected group. Any virtual spectrum is almost free of noise and only a limited number of *m/z* values are assigned with the respective intensities (*y*-coordinates). During its comparison with an experimental spectrum, it may theoretically happen that all peaks (or a vast majority of them) from the virtual spectrum will have their counterparts. A typical fingerprint pattern is discovered in this way, which is crucial for the unambiguous identification of spectra acquired from specific samples. Besides generating virtual spectra using the

BIOSPEAN, they can also be constructed outside the system and subsequently imported.

It becomes clear from the literature survey that both reference strains and studied isolates of microorganisms have to be analyzed by MALDI-TOF MS under the same experimental conditions to minimize undesired variations [10]. This means that not only the culture growth and sample preparation but also the instrumental setup for spectra acquisition should be optimized and standardized. Nevertheless, different culture conditions including media composition and temperatures of incubation may not necessarily affect obtaining good identification results [18]. For example, possible differences between the mass spectra from young and mature colonies can preferably be handled by adding two different reference spectra into the database: one acquired from young colonies and the other from mature colonies [19]. The sample preparation step is crucial. For IC/IS MS, prewashed spores are mixed on the MALDI target with an optimized matrix solution. To achieve reliable identification results, it is also necessary to use an updated commercial spectral database or a custom database engineered with the highest care. For fungi, it has been demonstrated that less common or uncommon species have to be included in the database otherwise low identification rates can be expected for real samples such as clinical mold isolates [18].

Another interesting feature of the BIOSPEAN resides in supporting a team work. Every user has to be registered for receiving a personal account. Data introduced through this account are password-protected and thus available only after a successful login (and of course to the system administrator). Nevertheless, this is not always desirable. In

the case of working on a project, where there are more participating researchers, it is desirable to share data. For doing that, “friends” may be selected from the group of registered users and granted with accessibility rights. Such an access is provided only to a given friend and it is usually restricted to a chosen group of spectra. When there is an intention to allow more groups of spectra accessible at the same time, the access to every group has to be adjusted individually. This strategy has been adopted for security reasons. On the other hand, any researcher may join more research teams and share their spectra.

## Technology

The BIOSPEAN has been developed as a web-based application and thus it benefits from a common web technology. Web applications are popular due to the convenience of using a web browser as a client. The main advantage of this solution resides in a centralized maintenance and updating, there are also easy requirements for any user (just a functional web browser). The computational kernel and logic were made in the language PHP, the user interface is a combination of HTML and JavaScript. Ajax was used as a web-development technology, MySQL was chosen for data storage because it is currently the most popular and widely used open source database technology [17]. All these technologies ensure that the application can be operated on a majority of web servers. Nevertheless, because of running background tasks (input analysis of spectra), the triggering mechanism requires running of the server using an operational system, which is compatible with Unix systems.

## Conclusions

The BIOSPEAN can be regarded as a novel solution for IC/IS MS and to our knowledge there is no comparable equivalent available in the form of a web application. It differs from the established commercial software for microbial biotyping in its concept and focus. The advantages of the BIOSPEAN include an easy-to-use graphical user interface, which is compatible with web browsers. Both software efficiency and reliability of the output information have already been demonstrated [16]. The application allows creating of a database for storage of MALDI-MS spectral data. Another important feature resides in the ability of virtual spectra construction. In consequence, this utility facilitates building of reference spectra for identification of specific samples. The BIOSPEAN also supports data sharing among users.

Thanks to the web technology used, datasets produced by a research team can be processed on more computers at once regardless the places from which users get their access. The software offers a controlled sharing and storage of data at a single place. Last but not least, the BIOSPEAN can be used with negligible operational costs. Due to the initial purchase costs, the common commercial software solutions become economical only when large data volumes are processed in a daily routine. But this can be done using the BIOSPEAN as well. Tests with an IC/IS MS identification of fungal microorganisms have shown its ability to work with thousands of spectra at the gigabyte data volume level. The BIOSPEAN is currently run in a test mode and after finding a long-term solution with respect to the technical support it will become available for a large public use. Then a further development of the software will also be facilitated, which is expected to include the addition of some statistical tools.

## Acknowledgements

This work was supported by OP RD&I grant no. ED0007/01/01 (Centre of the Region Haná for Biotechnological and Agricultural Research) and grant no. 7AMB12AT018 (student and academic staff mobilities within the Austria-Czech

Republic exchange program AKTION) from the Ministry of Education Youth and Sports, Czech Republic. We thank Dr. Martina Marchetti-Deschmann from the Vienna University of Technology, Austria, and Jana Chalupová from the Faculty of Science, Palacký University in Olomouc, for providing us with their data from IC/IS MALDI-TOF MS of fungi as a testing material.

## References

1. Yates JR (2001) Mass spectrometry in biology. In: Encyclopedia of Life Sciences. John Wiley & Sons Ltd, Chichester.
2. Griffiths WJ, Wang Y (2009) Mass spectrometry: from proteomics to metabolomics and lipidomics. Chem Soc Rev 38: 1882-1896.
3. Carbonnelle E, Mesquita C, Bille E, Day N, Dauphin B, et al. (2011) MALDI-TOF mass spectrometry tools for bacterial identification in clinical microbiology laboratory. Clin Biochem 44: 104-109.
4. Strohal M, Kavan D, Novák P, Volný M, Havlíček V (2010) mMass 3: a cross-platform software environment for precise analysis of mass spectrometric data. Anal Chem 82: 4648-4651.
5. Bauer C, Cramer R, Schuchhardt J (2011) Evaluation of peak-picking algorithms for protein mass spectrometry. Methods Mol Biol 696: 341-352.
6. Demirev PA, Fenselau C (2008) Mass spectrometry for rapid characterization of microorganisms. Annu Rev Anal Chem (Palo Alto Calif) 1: 71-93.
7. Dieckmann R, Erhard M, Kallow W, Saueremann S (2004) Method for the identification of microorganisms by mass spectrometry. EU patent EP1437673 A1.
8. Marko DC, Saffert RT, Cunningham SA, Hyman J, Walsh J, et al. (2012) Evaluation of the BrukerBiotyper and Vitek MS matrix-assisted laser desorption/ionization-time-of-flight mass spectrometry systems for identification of nonfermenting Gram-negative bacilli isolated from cultures from cystic fibrosis patients. J Clin Microbiol 50: 2034-2039.
9. Köhling HL, Bittner A, Müller KD, Buer J, Becker M, et al. (2012) Direct identification of bacteria in urine samples by matrix-assisted laser desorption/ionization time-of-flight mass spectrometry and relevance of defensins as interfering factors. J Med Microbiol 61: 339-344.
10. Posteraro B, De Carolis E, Vella A, Sanguinetti M (2013) MALDI-TOF mass spectrometry in the clinical mycology laboratory: identification of fungi and beyond. Expert Rev Proteomics 10: 151-164.
11. Karger A, Stock R, Ziller M, Elschner MC, Bettin B, et al. (2012) Rapid identification of *Burkholderia mallei* and *Burkholderia pseudomallei* by intact cell Matrix-assisted Laser Desorption/Ionisation mass spectrometric typing. BMC Microbiol 12: 229.
12. Firacative C, Trilles L, Meyer W (2012) MALDI-TOF MS enables the rapid identification of the major molecular types within the *Cryptococcus neoformans/C. gattii* species complex. PLoS One 7: e37566.
13. Bille E, Dauphin B, Leto J, Bournoux ME, Beretti JL, et al. (2012) MALDI-TOF MS Andromas strategy for the routine identification of bacteria, mycobacteria, yeasts, *Aspergillus* spp. and positive blood cultures. Clin Microbiol Infect 18: 1117-1125.
14. Keys CJ, Dare DJ, Sutton H, Wells G, Lunt M, et al. (2004) Compilation of a MALDI-TOF mass spectral database for the rapid screening and characterisation of bacteria implicated in human infectious diseases. Infect Genet Evol 4: 221-242.
15. Chalupová J, Sedlářová M, Helmelt M, Rehulka P, Marchetti-Deschmann M, et al. (2012) MALDI-based intact spore mass spectrometry of downy and powdery mildews. J Mass Spectrom 47: 978-986.
16. Franc V, Rehulka P, Raus M, Stulík J, Novak J, et al. (2013) Elucidating heterogeneity of IgA1 hinge-region O-glycosylation by use of MALDI-TOF/TOF mass spectrometry: Role of cysteine alkylation during sample processing. J Proteomics 92: 299-312.
17. Gibas C, Jambeck P (2001) Developing Bioinformatics Computer Skills.
18. Lau AF, Drake SK, Calhoun LB, Henderson CM, Zelazny AM (2013) Development of a clinically comprehensive database and a simple procedure for identification of molds from solid media by matrix-assisted laser desorption/ionization-time of flight mass spectrometry. J Clin Microbiol 51: 828-834.
19. Alanio A, Beretti JL, Dauphin B, Mellado E, Quesne G, et al. (2011) Matrix-assisted laser desorption ionization time-of-flight mass spectrometry for fast and accurate identification of clinically relevant *Aspergillus* species. Clin Microbiol Infect 17: 750-755.