# Rapid development of proteomic applications with the AIBench framework

**Hugo López-Fernández[1], Miguel Reboiro-Jato[1*], Daniel Glez-Peña[1], José R. Méndez Reboredo[1], Hugo M. Santos[2], Ricardo J. Carreira[2], José L. Capelo-Martínez[3], Florentino Fdez-Riverola[1]**

[1]Escuela Superior de Ingeniería Informática, University of Vigo, Edificio Politécnico, Campus Universitario As Lagoas s/n, 32004 Ourense, Spain

[2]REQUIMTE, Departamento de Química, Facultade de Ciências e Tecnologia, Universidade Nova de Lisboa, Lisboa, Portugal

[3]Bioscope Group, Physical Chemistry Department, University of Vigo, Ourense, Spain

### Summary

In this paper we present two case studies of Proteomics applications development using the AIBench framework, a Java desktop application framework mainly focused in scientific software development. The applications presented in this work are Decision Peptide-Driven, for rapid and accurate protein quantification, and Bacterial Identification, for Tuberculosis biomarker search and diagnosis. Both tools work with mass spectrometry data, specifically with MALDI-TOF spectra, minimizing the time required to process and analyze the experimental data.

## 1　Introduction and Motivation

Since the 1980s, computational methods have been used for comparative analysis of genome data [1], leading to the Bioinformatics field expansion, where various branches such as Genomics, Proteomics, Computer-Aided Drug Design, Bio Data Bases & Data Mining, Molecular Phylogenetics, Microarray Informatics or System Biology emerged [2]. With the arrival of high-throughput 'omics' data, a new dimension was added to data-driven comparative research [1], where manual analysis became unsuitable, requiring the use of computational methods. Nowadays, due to the proliferation of high-throughput techniques, the need for Bioinformatics tools is increasing, being necessary even in small studies.

In this context, the rapid development of successful feature-rich applications containing advanced functionalities in the field of biomedical and clinical research remains a major demand for smaller institutions due to the lack of both human and financial resources [3]. This fact gets worse if we consider the software development effort required to deliver highly specialized applications usually demanding sophisticated user interfaces. Moreover, developing applications in an interdisciplinary and applied research context also presents a large number of particular requisites ranging from computational requirements to usability. Specific issues include (*i*) sharing of heterogeneous data, (*ii*) integrating third party or previously developed

---

*To whom correspondence should be addressed. Email: mrjato@uvigo.es

algorithms, (*iii*) cross-platform compatibility, (*iv*) ability to repeat workflows but changing a few parameters or input data, (*v*) extensive use of logging messages to monitor the progress of long processes, (*vi*) establishing values for a high and variable number of parameters before running experiments and (*vii*) taking the maximum advantage of multi-threading capabilities in high-demanding tasks, among others.

Considering these requirements, we have developed AIBench [3], an open-source Java application framework for rapid development of translational software in Biomedicine. In this work, we describe the design of two new AIBench-based applications in the field of Proteomics: Decision Peptide-Driven and Bacterial Identification. Both tools analyze MALDI-TOF spectra in order to accurately quantify the amount of protein in a given sample and to discriminate between Tuberculosis strains, respectively.

The paper is structured as follows. Section 2 contains an overview of the AIBench architecture and its key design concepts. Section 3 identifies the scope of the applications developed and presents the shared components. Sections 4 and 5 describe the main functionalities and AIBench components of the Decision Peptide Driven and Bacterial Identification applications, respectively. Finally, Section 5 includes the conclusions and future work.

## 2   The AIBench Framework

The AIBench platform was particularly conceived to facilitate the development of a wide range of research applications based on general input-processing-output cycles where the framework acts as the glue between each executed task. In order to provide the basis for supporting rapid application development, AIBench manages three key concepts that are present in every AIBench application: operations, data-types and views. The developer only needs to concentrate on how to divide and structure the problem specific code into objects of these three entities.

From an architectonical perspective, AIBench is structured in several layers, as it is shown in Figure 1. The AIBench framework runs over a plug-in engine able to define a straightforward reusable component model where both the framework native components and the application-specific functionalities are divided and packaged into plug-ins. AIBench plug-ins are isolated by default, increasing the modularity and ensuring that accidental coupling is not introduced, but they can also interact by establishing dependencies or extension points. A dependency between plug-ins allows one plug-in to require other plug-ins to be present at runtime and to be allowed to access their classes and/or resources. An extension point declares a place where some plug-in can be extended by another plug-in (extension), usually providing a specific interface implementation.

The Core layer contains two native plug-ins: the Core and the Workbench. The AIBench Core detects and registers the application-specific operations, executes them upon request, keeps the results in the Clipboard structure and stores the session workflow in the History. The graphical user interface aspects are implemented in the Workbench plug-in, which creates the main application window, composes a menu bar with all the implemented operations, generates input dialogs when some operation is requested for execution, instantiates the registered results viewers, etc. All additional services bundled with AIBench belong to the Services layer and are also implemented via independent plug-ins that can be easily removed to meet the application
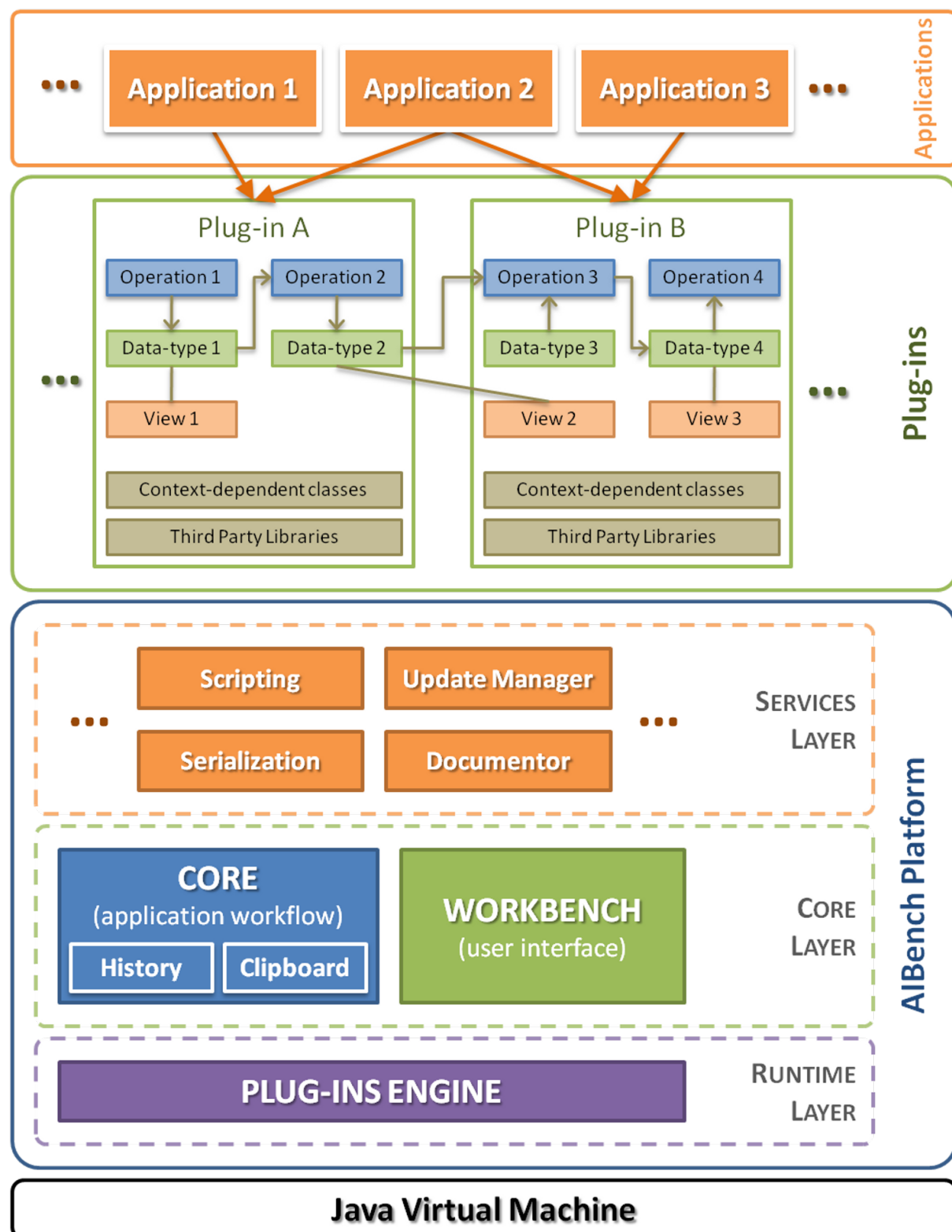
**Figure 1: AIBench framework architecture.**

specific needs. The Core and Services layers are maintained by the AIBench team and constitute all the code built-in and distributed with the framework, being the starting point of every development.

The Plugin and Application layers are placed on the top of the architecture and contains the application specific code (operations, data-types and views) provided by applications developers (AIBench users). In this sense, when an applications developer starts using the framework, there are no operations, data-types or views available, because these components are problem specific items. However, operations, data-types and views can (and should) be shared among applications related to the same area, especially when they are developed inside the same team. These higher-level components, along with other third-party libraries are also packaged in one or more plug-ins. Finally, from the most abstract point of view, an AIBench application can be seen as a collection of operations, data-types and views, reusable in more than one final application.

## 3   Building Proteomics Applications

Mass spectrometry (MS) is an analytical technique that generates mass spectra plotting the relative abundance of molecules versus their mass-to-charge ratio (m/z) [4]. This technique is an important tool in Proteomics for biomarker discovery, identification and quantification, as it allows the detection of small molecules (such as peptides) and the measurement of their intensities. Due to the low analysis time and the high data throughput of MS, it is very useful for fast and accurate diagnosing and monitoring different disease states [5]. On the other hand, the volume of data generated in MS experiments makes manual analysis unsuitable, requiring the use of Bioinformatic applications.

In this context, we have used the AIBench framework for the development of two MS data analysis applications with specific objectives, which will be presented in the following sections. These applications are intended to work with MALDI-TOF data, a soft ionization technique used in MS. As the data source for both applications is the same, a first plug-in named 'Peptide Analysis' has been developed containing common data-types, views and operations. Figure 2 shows the main classes of this plug-in.

The SpectrumData class represents a MS spectrum composed by a set of peaks each with mass-to-charge ratio and intensity. ExperimentalData contains a set of SpectrumData coming from the same sample or data analysis. These three classes are specializations of three more generic classes, which can be used to work with other types of experimental data. The LoadSpectrumData and LoadExperimentalData operations can be used to load spectra from CSV files and, although it is not represented in the Figure 2, the ExportToCSV operation is a generic operation to store different data-types into CSV files. Finally, the views in the Figure 2 present the main data-types as charts and tables.

The following sections describe the case studies of the two Proteomics applications developed with the AIBench framework and based on the 'Peptide Analysis' plug-in previously introduced.
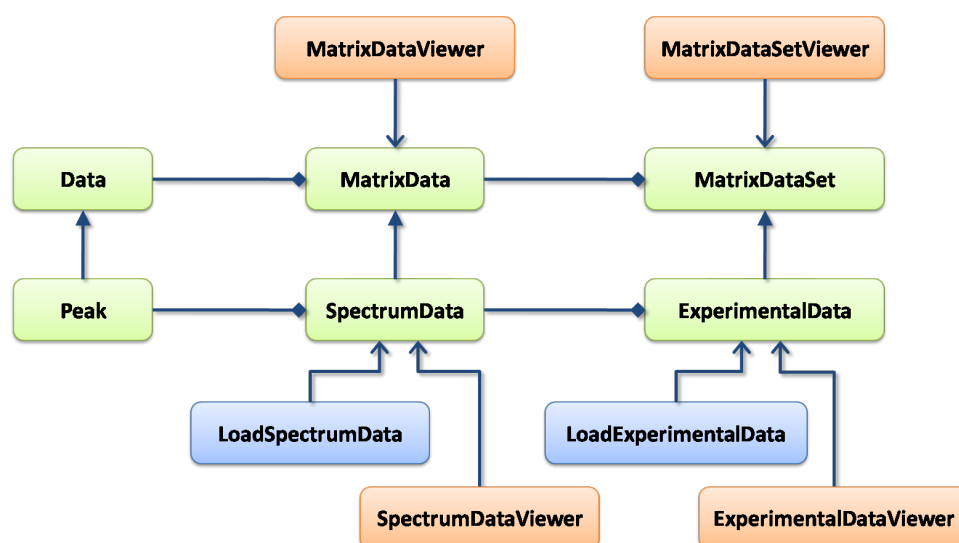
**Figure 2: Main data-types (green), operations (blue) and views (orange) of the 'Peptide Analysis' plug-in.**

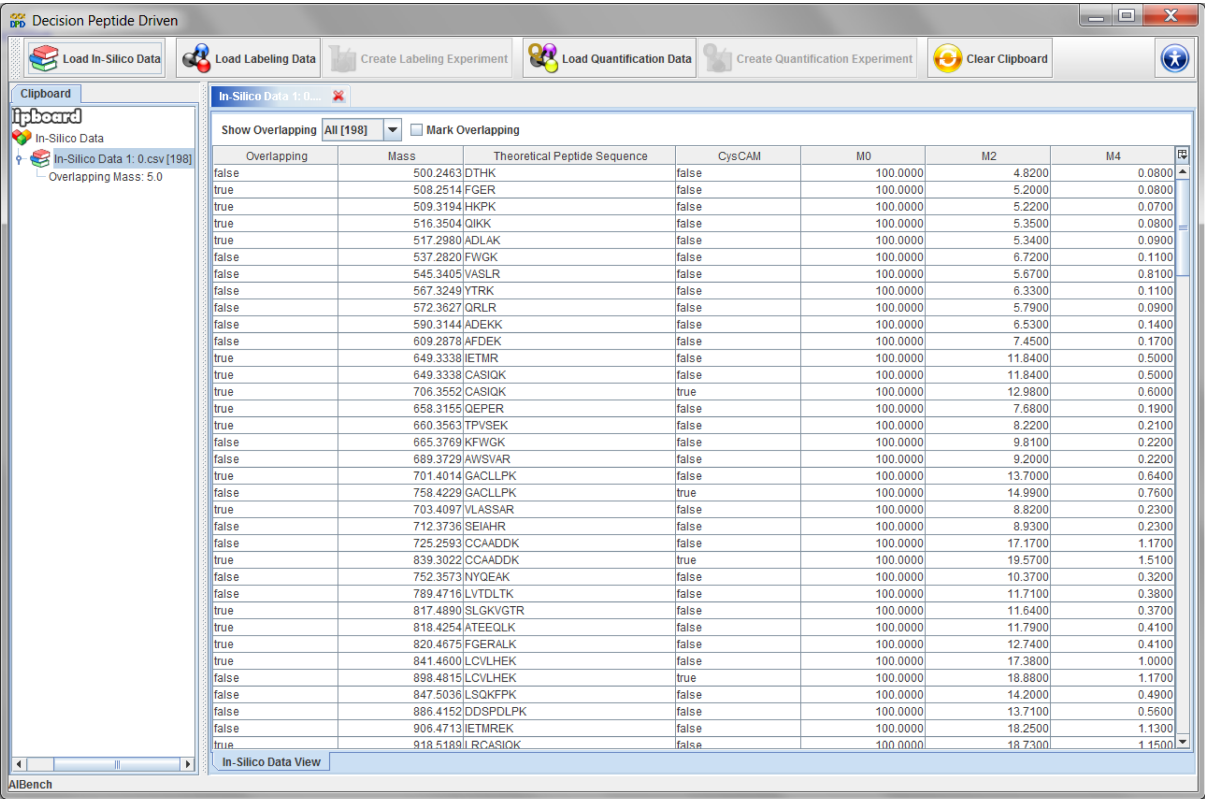## 4 Case Study 1: Decision Peptide-Driven

The identification and quantification of the protein content of biological samples plays a crucial role in biological and biomedical research [6]. The Decision Peptide-Driven tool[1] (DPD) [7] implements a software application for assisting the user in a protocol [8] for accurate protein quantification by using MALDI-TOF mass spectrometry. Main application views showing data and experiment information are shown in Figures 3 and 4, respectively.

Using the DPD software, user can compare the MALDI results of the direct and inverse $^{18}$O-labeling experiments and quickly identify those peptides with paralleled loses in different sets of a typical proteomic workflow. Those peptides are used for subsequent accurate protein quantification. The interpretation of the MALDI data from direct and inverse labeling experiments is time-consuming requiring a significant amount of time to do all comparisons manually.
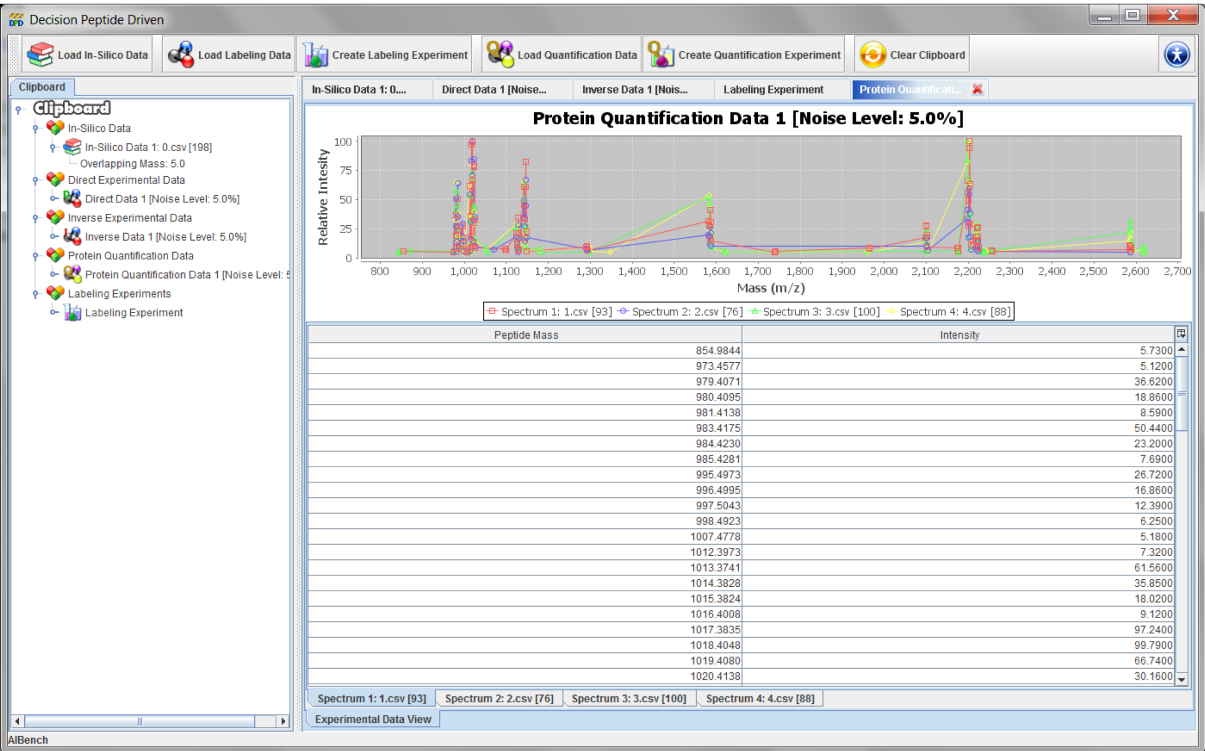
The DPD software shortens and simplifies the search for the peptides that must be used for quantification from a week to just some minutes. To do so, it takes as input several MALDI spectra and aids the researcher in an automatic mode (*i*) to compare data from direct and inverse $^{18}$O-labeling experiments, calculating the corresponding ratios to determine those peptides with paralleled losses throughout different sets of experiments; and (*ii*) allow to use those peptides as internal standards for subsequent accurate protein quantification using $^{18}$O-labeling.

When working with the DPD software, users can follow two main workflows: labeling and quantification. Figure 5 depicts the data-types (green) and operations (blue) involved in the labeling workflow and their relationships. Along with direct and reverse data, *in silico* data is needed to perform the peptide detection in the spectra. Once the data has been loaded, the reproducible peptides can be obtained by executing a labeling experiment. This is a macro-operation composed by a series of operations: (*i*) detect peptides, indentifies the peptides present in the direct and inverse data using the *in silico* data as reference, (*ii*) filter peptides, removes the low-quality peptides, (*iii*) intersec peptides, joins the direct and inverse peptides, and (*iv*) select reproducible, evaluates the quality of the final peptides. The final output of the labeling work-
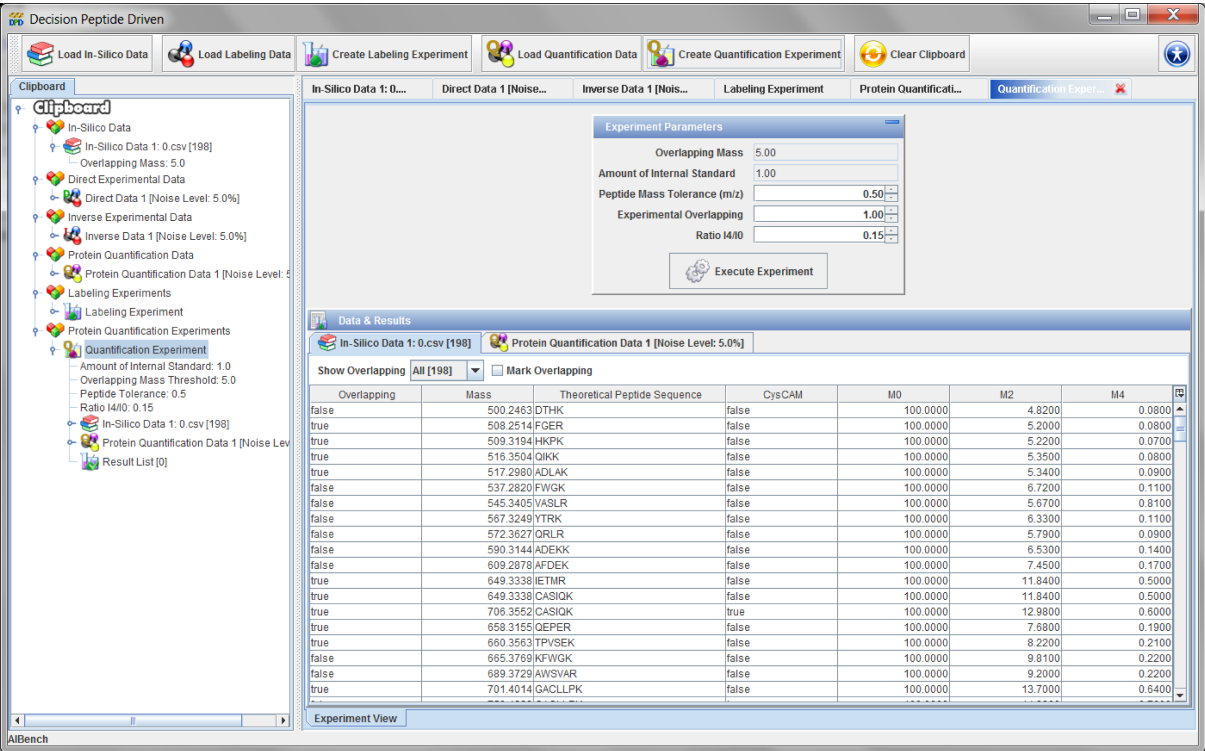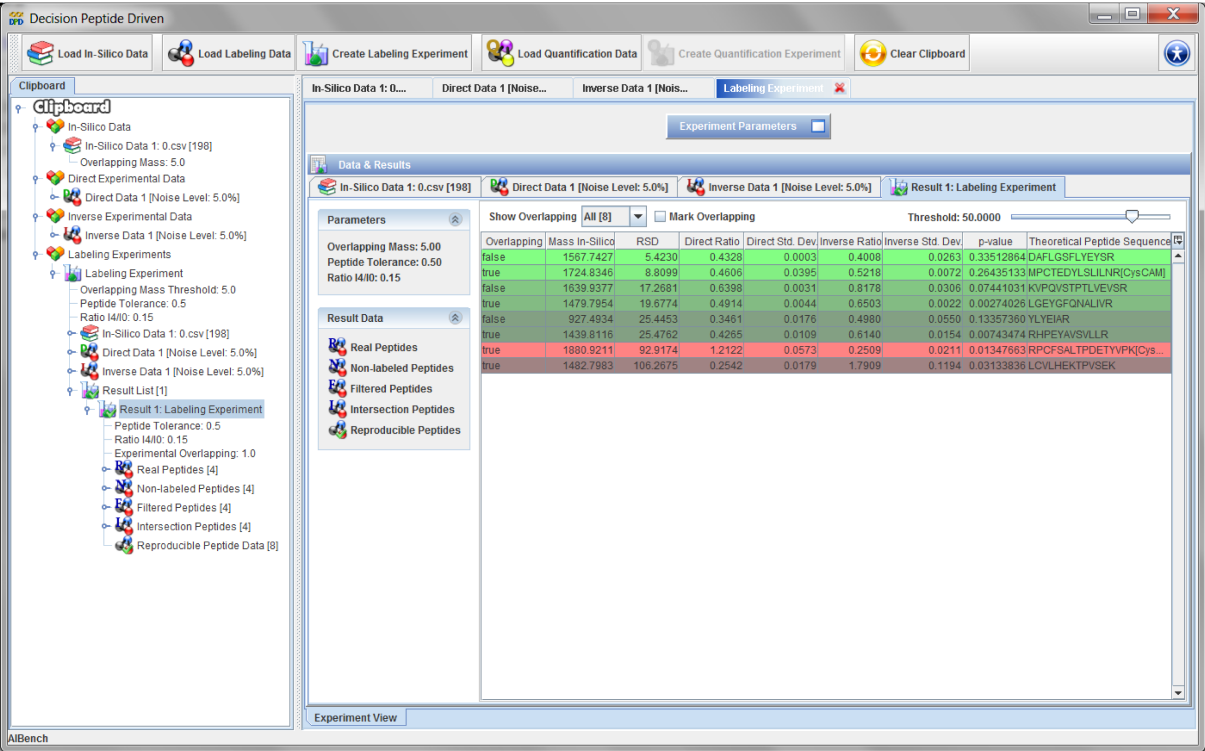
---

[1]http://sing.ei.uvigo.es/DPD

**(a)**



**(b)**

**Figure 3: Snapshots of (*a*) the *in silico* data view and (*b*) the experimental data view of the DPD software.**
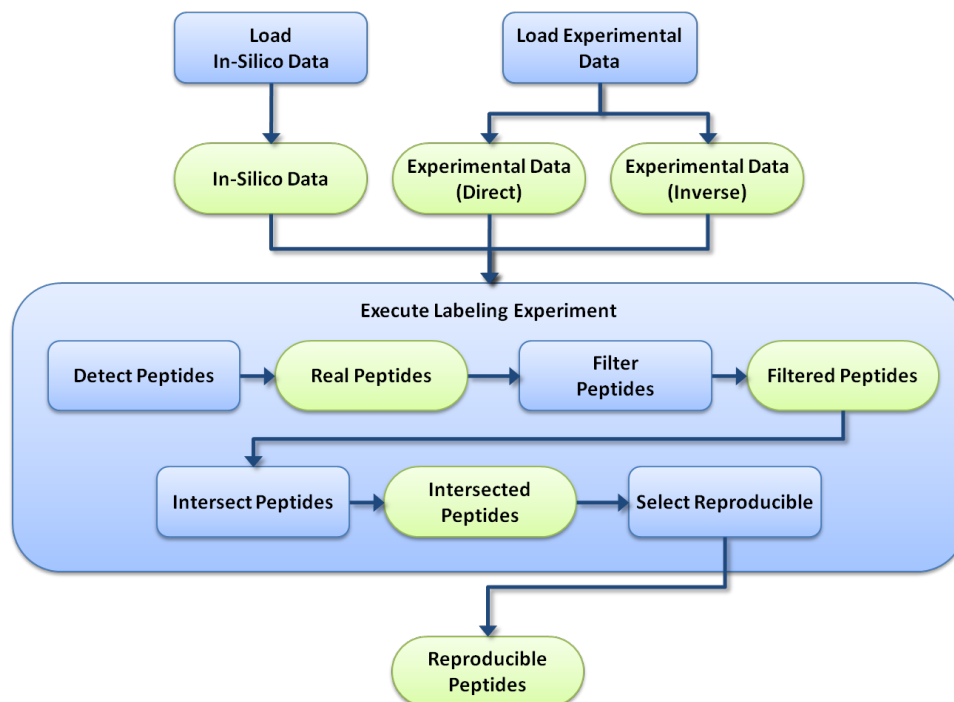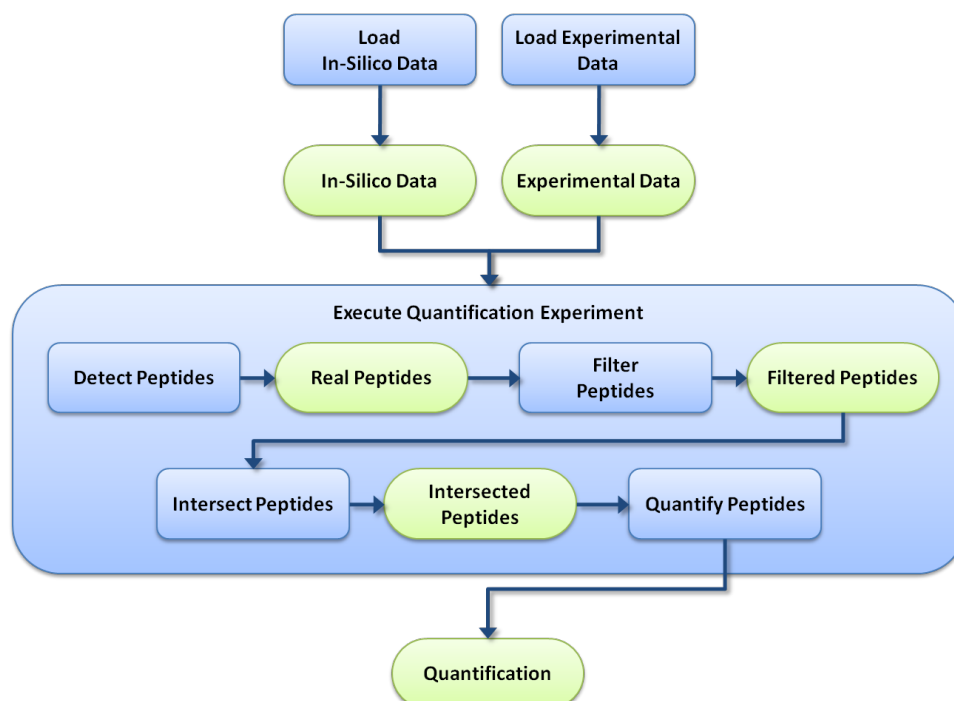
**(a)**



**(b)**

**Figure 4:** Snapshots of (*a*) the experiment view and (*b*) the experiment results view of the DPD software.

flow is a list of reproducible peptide with a quality value, along with the data generated in the intermediate steps of the experiment.

**Figure 5: DPD labeling experiment workflow.**

The quantification workflow allows the protein quantification in a MALDI-TOF sample. As can be seen in Figure 6, this workflow is very similar to the labeling workflow. The main difference is that, in this case, only one experimenal data is needed. Additionally, the last operation of the quantification experiment is a quantification step where the amount of protein is calculated.

**Figure 6: DPD quantification experiment workflow.**

Both workflows described above are sequential workflows, as they implement two concrete protocols for MS data analysis specifically designed for the proposed problems. This type of protocols suits perfectly to the input-process-output model on which is based AIBench, thus we could take advantage of the framework capabilities in order to minimize the development time.

The AIBench components used to build the DPD application are distributed in two plug-ins: the 'Peptide Analysis' plug-in, described in the previous section, and the 'Decision Peptide-Driven' plug-in, containing all operations, data-types and views of the software. The second plug-in depends on the first one. The list of components included in the 'Decision Peptide-Driven' plug-in are summarized in Table 1.

Table 1: AIBench components of the DPD software.

| Component | Description |
|---|---|
| *Operations* | |
| Create Labeling Experiment | Creates a new labeling experiment. |
| Create Quantification Experiment | Creates a new quantification experiment. |
| Execute Experiment | Executes a labeling or quantification experiment. |
| Export Result | Exports an experiment result as a log file. |
| Load In Silico Data | Loads *in silico* peptide information from a csv file. |
| Load Labeling Data | Loads experimental spectra data for labeling experiments from csv files. |
| Load Quantification Data | Loads experimental spectra data for quantification experiments from csv files. |
| *Data-types* | |
| InSilicoData | *In silico* peptide information. |
| DirectExperimentalData | Experimental spectra data to be used as direct data in labeling experiments. |
| InverseExperimentalData | Experimental spectra data to be used as inverse data in labeling experiments. |
| LabelingExperiment | A labeling experiment. Contains direct and inverse data and the parameter configuration. |
| QuantificationExperimentalData | Experimental spectra data to be used in quantification experiments. |
| QuantificationExperiment | A quantification experiment. Contains quantification data and the parameter configuration. |
| IExperiment | Common interface for labeling and quantification experiments. |
| IResult | Common interface for experiment results. |
| Storeable | Items that can be exported to log files. |
| *Views* | |
| Experiment View | Allows the user to execute experiments and view results. |
| Experimental Data View | Shows chart and table representation of the experimental spectra data. |
| In Silico Data View | Shows table representation of the *in silico* data. |

# 5 Case Study 2: Bacterial Identification

Tuberculosis (TB) is an infectious disease and the world's leading cause of death. In fact, the World Health Organization estimates that one third of the world's population is infected with Mycobacterium tuberculosis (MTB) [9]. Moreover, the usual methodologies for TB diagnosis can take up to 3 months for the identification of the correct infectious agent or only 1 week with molecular methods like PCR, but these last methods are expensive and are not affordable for many laboratories, especially in undeveloped countries. Recently, Hettick *et al.* reported the discrimination of intact mycobacteria at the strain level using MALDI-TOF-MS and a bio-statistical analysis [10]. However, for security reasons the handling of intact TB, MDR-TB and XDR-TB can only be done in laboratories with a biosafety level of 3, and thus the analysis of intact MTB by mass spectrometry is not possible.

Bacterial Identification tool[2] (BI) is a software application that allows rapid bacterial identification through peptide mass mapping obtained after analysis by MALDI-MS-based approaches.

The main functionalities of BI tool allow the researcher to (*i*) load and visualize spectra data containing peaks of each bacterial strain, (*ii*) identify a peptide fingerprint of each strain, containing all common peaks of all replicates, (*iii*) fingerprint comparison to obtain peaks which are unique of a given strain, being potential biomarkers, (*iv*) train and test a Support Vector Machines (SVM) classifier using the most predictive peptides and (*v*) perform and visualize a hierarchical clustering of the strains with those predictive peptides. Main views of the application are shown in Figures 7 and 8.
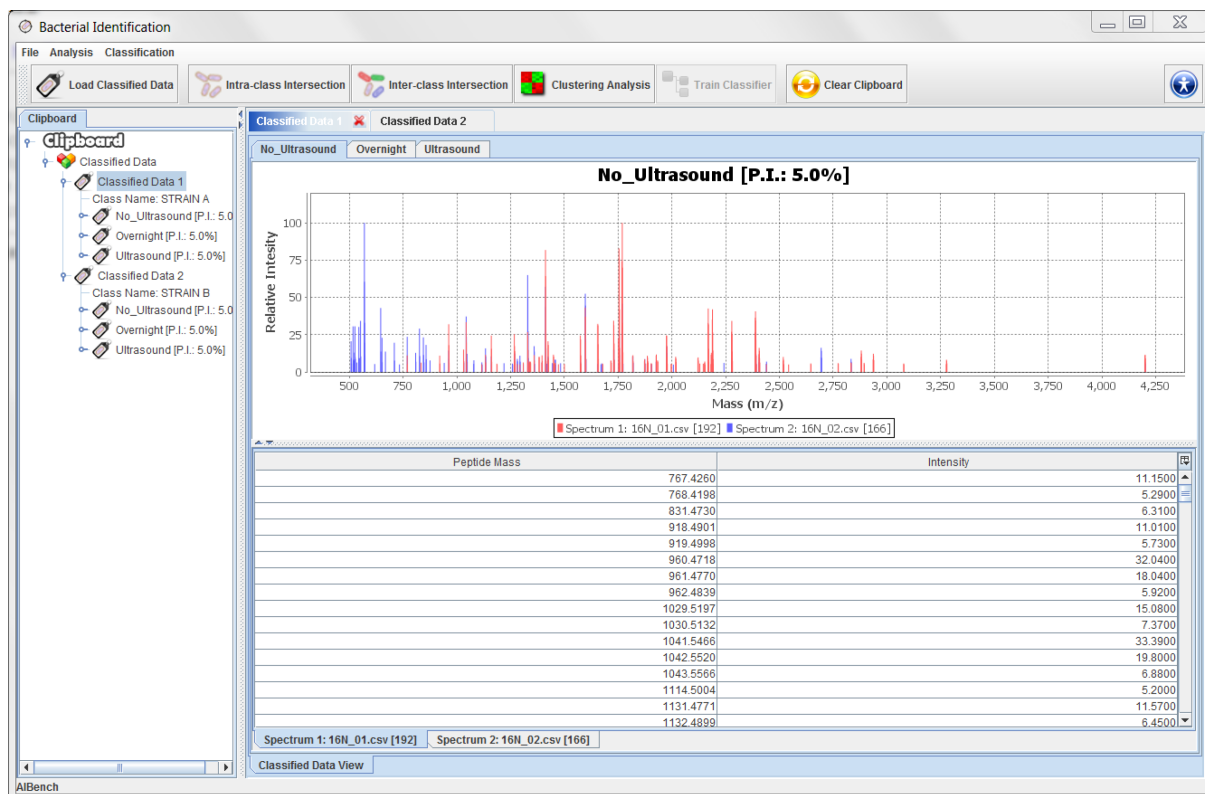
Figure 9 shows the internal workflow of the BI tool. Besides the data loading, this workflow can be divided into three sub-workflows, according to the functionalities they cover. The 'Strains Analysis' sub-workflow includes the 'intra' and 'inter' class analysis, useful for identifying the peptide fingerprint of TB strains and the discriminant peaks among different strains. The 'Classification' sub-workflow allows the user to train a SVM classifier using TB data previously classified and the correspondent list of discriminant masses. Once the classifier is trained, it can be used to classify new TB data. Finally, in the 'Clustering' sub-workflow, user can perform a hierarchical clustering, creating a heat map to study the relationship between the analyzed strains and, additionally, between their peptides.

In this case, we do not have a sequential workflow but, as noted above, the main workflow is divided into three different sub-workflows according to the analysis to perform. As the data loaded can be used in different analysis, the AIBench capability to keep data on the Clipboard so you can reuse it in different operations was very useful when developing this application.
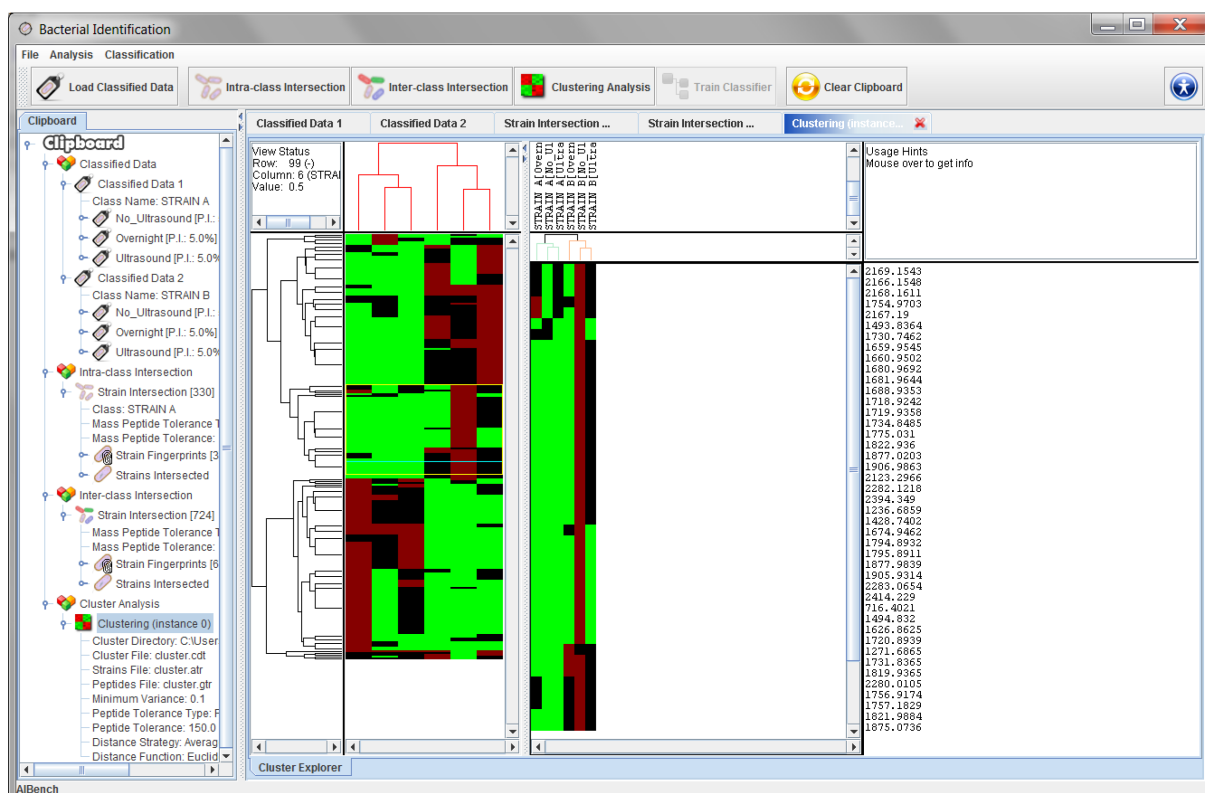
The AIBench components used to build the BI application are distributed in two plug-ins: the 'Peptide Analysis' plug-in and the 'Bacterial Identification' plug-in, containing all operations, data-types and views of the software. As in the previous case, the second plug-in depends on the first one. The list of components included in the 'Bacterial Identification' plug-in are summarized in Table 2.
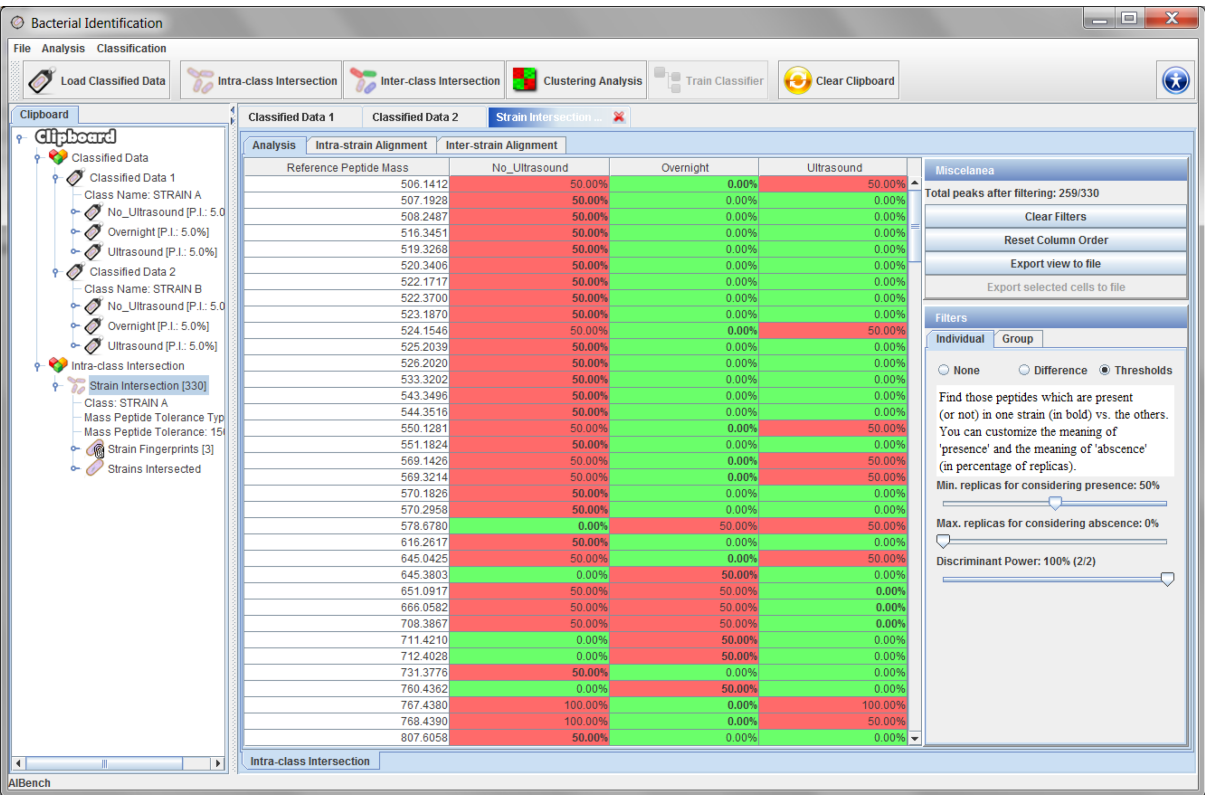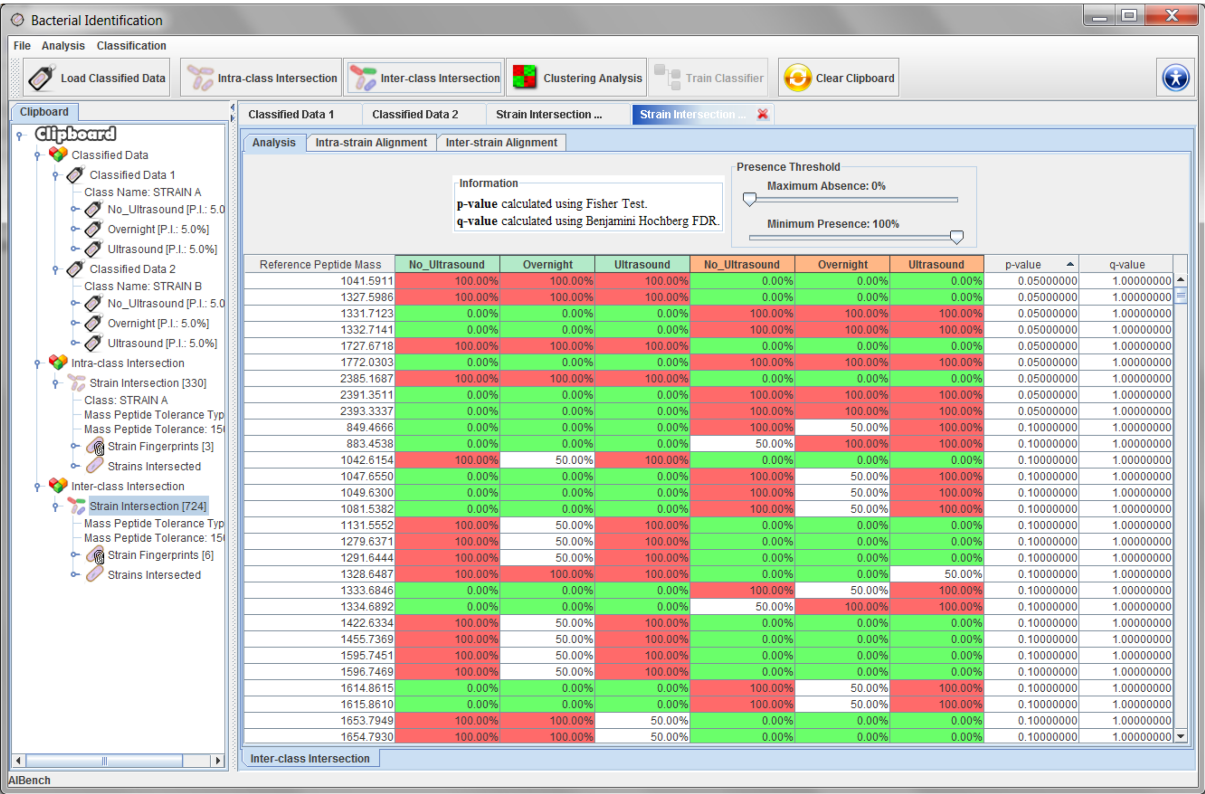
---

[2]http://sing.ei.uvigo.es/BI

**(a)**



**(b)**

**Figure 7:** Snapshots of (*a*) the data view and (*b*) the clustering view of the BI Software.

**(a)**



**(b)**

**Figure 8:** Snapshots of (*a*) the intra-class analysis view and (*b*) the inter-class analysis view of the BI Software.

Table 2: AIBench components of the Bacterial Identification software.

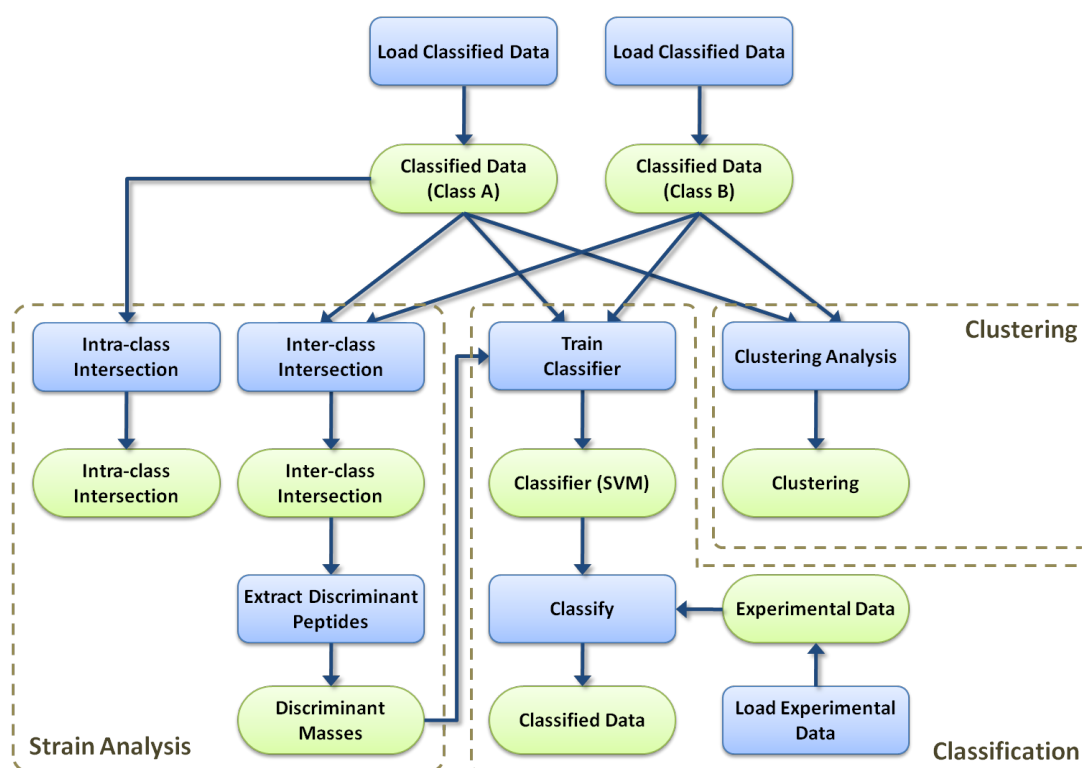| Component | Description |
| --- | --- |
| *Operations* | |
| Clustering Analysis | Performs a hierarchical clusternig analysis. |
| Create Discriminant Peptide List | Creates a new discriminant peptide list. |
| Inter-class Intersection | Intersect several strain data with different class and generates a peptide occurrence matrix. |
| Intersect Discriminant Peptide Lists | Intersects two or more discriminant peptide lists. |
| Intra-class Intersection | Intersects several strain data with the same class and generates a peptide occurrence matrix. |
| Load Classified Data | Loads several strains data of the same class. Strains data must be in separated directories. |
| Load Clustering | Loads clustering analysis results. |
| Load Discriminant Peptide List | Loads a stored list of discriminant peptides. |
| Load Trained Classifier | Loads a trained classifier. |
| Merge Discriminant Peptide Lists | Merges two or more discriminant peptide lists. |
| Save Discriminant Peptide List | Stores a list of discriminant peptides. |
| Save Trained Classifier | Stores a trained classifier. |
| Single Inter-class Intersection | Intersects several strain data with different class and generates a peptide occurrence matrix. |
| Single Intra-class Intersection | Intersects several strain data with the same class and generates a peptide occurrence matrix. |
| Substract Discriminant Peptide Lists | Subtract two or more discriminant peptide lists to a reference peptide list. |
| Train Classifier | Trains a new SVM-based strain classifier. |
| *Data-types* | |
| ClassifiedData | Strain data associated with a class. |
| Clustering | Clustering data. |
| DiscriminantPeptideList | List of discriminant peptides. |
| InterClassIntersection | Peptide occurrence matrix coming from an inter–class intersection. |
| IntraClassIntersection | Peptide occurrence matrix coming from an intra–class intersection. |
| TrainedClassifier | Trained strain classifier. |
| PeptideFingerprints | List of characteristic peptides of a class. |
| *Views* | |
| Classified Data View | Shows chart and table representation of the strains data. |
| Cluster Explorer | Shows a heat map and hierarchical view of clustering data. |
| Data Classification | Allows the user to train a classifier and to classify data with it. |
| Inter-class Intersection | Allows the user to analyze a peptide occurrence matrix coming from an inter-class intersection. |
| Intra-class Intersection | Allows the user to analyze a peptide occurrence matrix coming from an inter-class intersection. |

**Figure 9: Bacterial Identification workflow.**

# 6    Conclusions

This paper has presented two new applications, Decision Peptide-Driven and Bacterial Identification, in the field of Proteomics, showing the suitability of our AIBench framework to fast-develop user friendly applications in Biomedicine. The flexibility of AIBench has become evident as it was easily adapted to the different workflows of both applications. Additionally, both applications share components via common plug-ins, one of the most powerful aspects of AIBench.

The future work is divided in two objectives. By one hand, the two presented applications will be improved in their future versions. Decision Peptide-Driven will be integrated with Swissprot in order to download *in silico* data from the desired proteins and Bacterial Identification will include more classifier algorithms. In the other hand, the AIBench framework is in continuous development, including new features and bug fixes frequently.

## Acknowledgements

# References

[1] P. Hogeweg. The roots of bioinformatics in theoretical biology. PLoS Computational Biology, 7(3):e1002021. 2011.

[2] A.S. Nair. Computational biology & bioinformatics: a gentle overview. Communications of the Computer Society of India. 2007.

[3] D. Glez-Peña, M. Reboiro-Jato, P. Maia, M. Rocha, F. Díaz and F. Fdez-Riverola. A rapid application development framework for translational research in biomedicine. Computer Methods and Programs in Biomedicine, 98(2):191-203, 2010.

[4] O.D. Sparkman. Mass spectrometry desk reference. Global View Publishing. 2000.

[5] P. Findeisen and M. Neumaier. Mass spectrometry based proteomics profiling as diagnostic tool in oncology: current status and future perspective. Clinical Chemistry and Laboratory Medicine, 47(6):666-684, 2009.

[6] X.J. Li, E.C. Yi, C.J. Kemp, H. Zhang and R. Aebersold. A software suite for the generation and comparison of peptide arrays from sets of data collected by liquid chromatography-mass spectrometry. Molecular & Cell Proteomics, 4(9):1328-1340, 2005.

[7] H.M. Santos, M. Reboiro-Jato, D. Glez-Peña, J.D. Nunes-Miranda, F. Fdez-Riverola, R. Carvallo and J.L. Capelo. Decision Peptide-Driven: a free software tool for accurate protein quantification using gel electrophoresis and matrix assisted laser desorption ionization time of flight mass spectrometry. Talanta, 82(4):1412-1420, 2010.

[8] H.M. Santos, D. Glez-Peña, M. Reboiro-Jato, F, Fdez-Riverola, M.S. Diniz, C. Lodeiro and J.L. Capelo-Martínez. A novel 18O inverse labeling-based workflow for accurate bottom-up mass spectrometry quantification of proteins separated by gel electrophoresis. Elecrophoresis, 31:3407-3419, 2010.

[9] D. Kirschner. Timebomb: the global epidemic of multidrug resistant tuberculosis. Nature Medicine, 7:1173-1174, 2001.

[10] J.M. Hettick, M.L. Kashon, J.E. Slaven, Y. Ma, J.P. Simpson, P.D. Siegel, G.N. Mazurek and D.N. Weissman. Discrimination of intact mycobacteria at the strain level: a combined MALDI-TOF MS and biostatistical analysis. Proteomics, 6(24):6416-6425, 2006.