

A Laboratory Information Management System (LIMS) for High-Throughput LC-MS Metabolomics-Based Biomarker Discovery



Alan M. Smith, Ph.D.¹, Yuerong Zhu, Ph.D.², Paul R. West, Ph.D.¹, April M. Weir, M.S.¹, Gabriela G. Cezar, DVM, Ph.D.¹
¹Stemina Biomarker Discovery, Inc., 504 S. Rosa Rd., Suite 150, Madison, WI 53719
²BioInfoRx, Inc., 1502 Beechwood Cir, Middleton, WI 53562

OVERVIEW

- Developed a LIMS system combining experiment, sample, LC-MS, feature analysis, and interpretation of data into an interconnected web-based user interface
- System designed to handle thousands of LC-MS samples from planning to final interpretation
- Open source framework allows rapid integration of new project specific data analysis tools
- System tested on the Stemina Developmental Toxicity Platform
- Complex analysis procedures were simplified into simple web-based interfaces allowing nearly complete data analysis automation

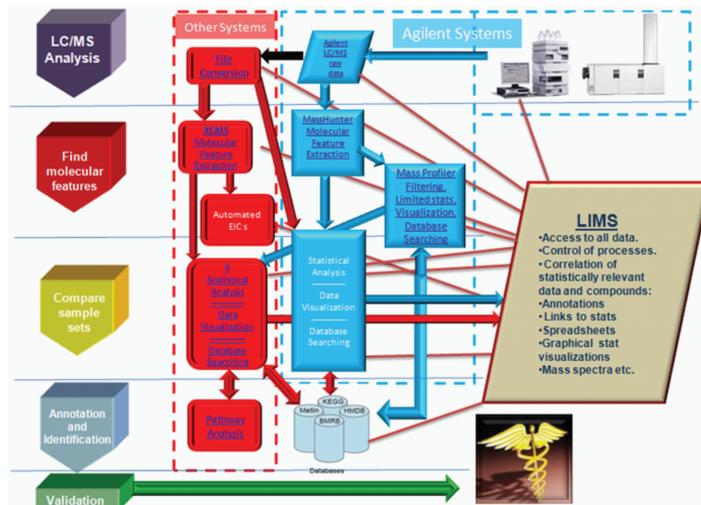


Figure 1. General overview of analysis pipeline combining proprietary and open source software managed by the SteminaLIMS.

INTRODUCTION

Large sample size multifactorial experiments executed under high throughput conditions require laboratory information management systems (LIMS) capable of simplifying an entire project pipeline (Figure 1). SteminaLIMS is a web-based laboratory management tool specifically designed to streamline the processes involved with metabolomics-based biomarker discovery (Figure 2). The SteminaLIMS manages the entire stream of textual, numerical, statistical and graphical data including project information, sample preparation and metadata, mass spectrometry data acquisition and analysis, file conversion, peak picking, statistical analysis, annotation, and small molecule confirmation by MS-MS. The system incorporates convenient user friendly tools that allow collaborative input, tracking, and dissemination of information. The data analysis pipeline uses both proprietary and open source data analysis tools to simplify the trek of small molecule biomarker discovery.

Project Tools

1. Projects
 2. Experiments
 3. Samples for MS Run - **New Version**
 4. MS Run Information
 5. MS Data Feature Extraction
 6. MS Data Profiler Filtering
 7. Statistical Analysis Summary
 8. Summary Table
- DevTox Project Progress
 - Culture Plate Management
 - Plot EICs
 - Plot EICs for non-DevTox
 - Upload m/z Data
 - Statistical Analysis
 - MS Data Interpretation
 - MS Data Validation
 - Define Project Fields
 - Experiment Sample Info

Inventory System

- **Stemina Laboratory Ordering System**
Submit ordering requests, process ordering, and check off received items.

Other Tools

- **Share files among group members**
- Stemina Small Business Server (external link)
- Store my files online
- Submit requests to group members
- Write daily notes
- Send e-mails to group members
- Post your Blogs
- **Update my profile** - change password, e-mail address

No.	Blog Subject	Comments
1	Using the LIMS system from home	1
2	Overall Project Progress Check List	0
3	Recent Project Progress by Ron	0

Figure 2. LIMS interface for Developmental Toxicity metabolomics platform. A large number of tools are available for the user to manage and analyze data.

METHODS

The SteminaLIMS is built by BioInfoRx based on its innovative information management technologies. Web programming technologies, such as LAMP (Linux, Apache, MySQL, and PHP) and AJAX (Asynchronous JavaScript and XML), are used to design the backend databases and the web-based interfaces. The system integrates data management with LC-MS analysis processes, providing improved work flow automation. Custom statistical analysis tools were created with open source statistical software R and Bioconductor. These tools allow users to select and preprocess mass feature tables, bin mass features across LC-MS samples, perform differential analysis of feature abundances and create summary tables to identify features of interest. Online mass spectral data analysis tools for EIC and spectra generation are based on the XCMS library in R.

RESULTS

Drug	MS Date	MS Step	Polarity	PF ID
Ignore	Ignore	Ignore	Ignore	Ignore
5-Fluorouracil	2008-07-18	B12	Neg	PF2~VPA2~2008-07-18~Pos
Ascorbic_Acid	2008-07-29	C14	Pos	PF4~Indomethacin~2008-07-29~Pos
Aspirin	2008-07-30			PF6~Hydroxyurea~2008-07-30~Pos
BT22 Treated	2008-08-01			PF9~Aspirin~2008-08-01~Pos
BT33 Treated	2008-08-07			PF10~VPA2~2008-08-07~Neg
BTSC12.1 Treated	2008-08-12			PF11~Indomethacin~2008-08-07~Neg
BTSC22 Treated	2008-08-13			PF12~Hydroxyurea~2008-08-07~Neg
BTSC33	2008-08-18			PF13~Aspirin~2008-08-07~Neg
	2008-08-19			PF14~Thalidomide~2008-08-12~Neg

Analysis Name: TestRun
 Image Size: Width: 16 (inches) - Height: 16 (inches), Font Size: 14, Resolution: 72

Abundance cutoff: 0 (number) Removes features under this abundance
 Missing Negatives: 10000 (number) Replace missing negative value
 Missing Positives: 50000 (number) Replace missing positive value
 NAallowableC: 0.5 (0-1.0) Threshold of missing data to remove LC-MS run
 NAallowableR: 1 (0-1.0) Threshold of missing data to remove feature bin
 Bin ppm1: 1e-5 NUMBER; the ppm cutoff for features under 175da
 Bin ppm2: 7e-6 NUMBER; the ppm cutoff for features 175-300da
 Bin ppm3: 5e-6 NUMBER; the ppm cutoff for features 300+ da
 RTcut: 12 (seconds) Retention time cutoff
 Adunfilter: 0 NUMBER; Remove abundance less than this value
 Massfilter: 0 NUMBER; Remove mass less than this value
 RTfilter: 0 NUMBER; Remove RT less than this value
 Abundcol: yes (text "yes" or "no"); # does the input file have an abundance column input YES or N
 Mean Center: True False
 SD scale: True False
 annocutoff: 15 NUMBER ppm; excludes features below X pvalue
 pvalcutoff: 1 (number 0-1.0); excludes features below X pvalue
 DBname: Choose:

Figure 3. Example of user interface to analyze outputs from Agilent's Mass profiler software. This data analysis procedure reads in Mass Profiler outputs, reorganizes the data, groups mass features across experiments, performs univariate and multivariate statistical analysis, creates project specific graphical outputs of the analysis, annotates the features using several in-house DB's, and generates easy to read data summary documents.

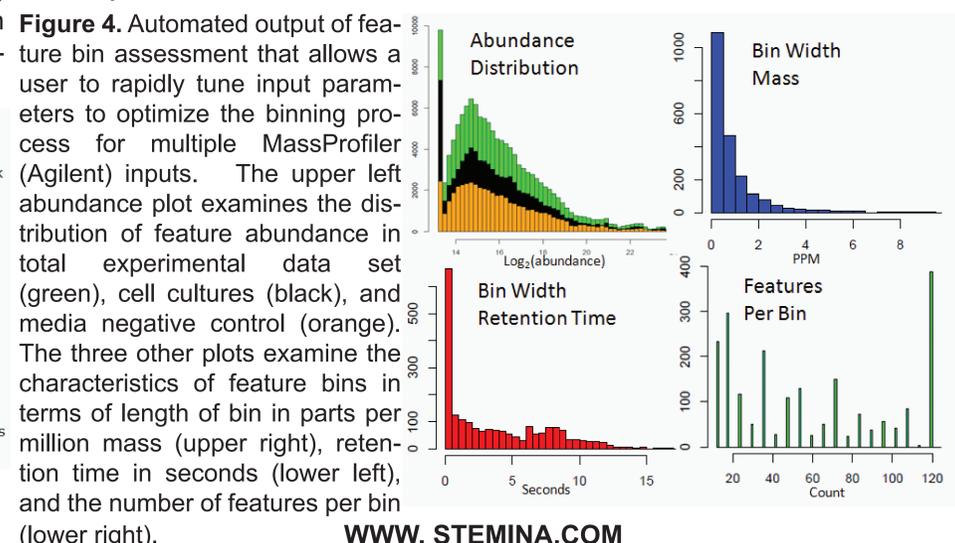


Figure 4. Automated output of feature bin assessment that allows a user to rapidly tune input parameters to optimize the binning process for multiple MassProfiler (Agilent) inputs. The upper left abundance plot examines the distribution of feature abundance in total experimental data set (green), cell cultures (black), and media negative control (orange). The three other plots examine the characteristics of feature bins in terms of length of bin in parts per million mass (upper right), retention time in seconds (lower left), and the number of features per bin (lower right).

WWW.STEMINA.COM

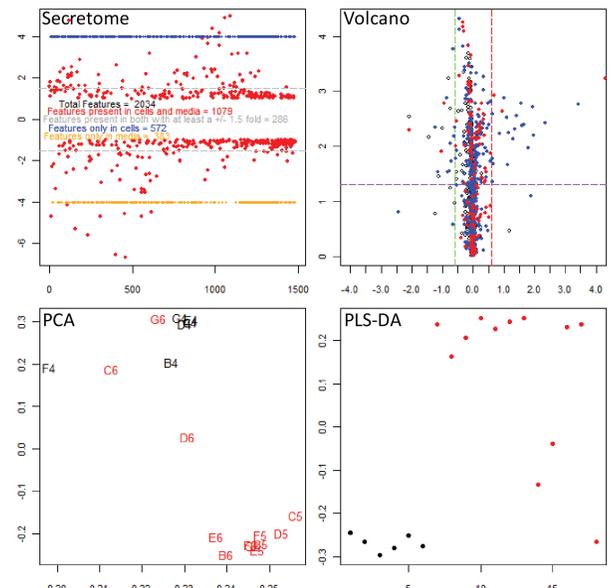


Figure 5. LIMS generated plots based on the results of summary information, univariate, and multivariate statistics. The top left plot evaluates the secretome by changes in abundances of features with respect to cell cultures (+ cells) and media negative controls (- cells). The top right volcano plot evaluates the response to treatment of features present in the cell culture. For both plots: features detected only in the presence of cells (blue), features present in both cells and media (red), features present only in the media (orange), features present in both cells and media, but not above the threshold lines (open black circles). The bottom left loadings plot of the samples contains the first and second principle components from NIPALS PCA (Black = Control, Red=Treated). The bottom right PLS-DA scores plot demonstrates separation of treated (red) and control (black) samples using the SIMPLS method. Information used in each of these plots is also provided in summary spreadsheets such as difference from media, p values, and VIP scores for each feature.

The bottom left loadings plot of the samples contains the first and second principle components from NIPALS PCA (Black = Control, Red=Treated). The bottom right PLS-DA scores plot demonstrates separation of treated (red) and control (black) samples using the SIMPLS method. Information used in each of these plots is also provided in summary spreadsheets such as difference from media, p values, and VIP scores for each feature.

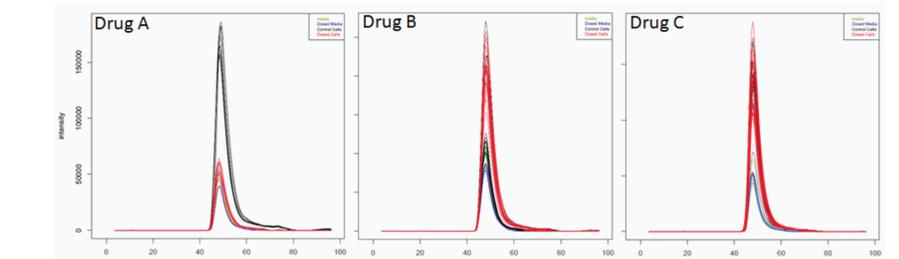


Figure 6. Smoothed extracted ion chromatogram (EIC) for a single mass feature generated by the LIMS (based on features that are secreted by human embryonic stem cells) exhibits statistically significant difference in abundance in response to specific classes of drugs, and were selected based on variable importance from a PLS-DA model.

CONCLUSIONS

- An open-source LIMS creates a working environment that can simplify a complex data analysis pipeline (Figure 1) into a simple user interface (Figure 2).
- Project specific data analysis procedures can be reduced to simple and easy to use web-based user interfaces with minimal effort (Figure 3).
- Tailed outputs for parameter optimization such as combining large numbers of MassProfiler outputs can be easily presented to system users who can evaluate the results (Figure 4).
- Project specific outputs and summary statistics are created and presented to users freeing bioinformaticists for other projects (Figure 5).
- The final results of statistical analysis can be evaluated by users for features that are acceptable for further confirmation by a mass spectrometrist (Figure 6).
- Project data was entered by multiple scientific specialists and passed successfully from experimental design stages to final MS-MS confirmation of differentially secreted small molecules using the SteminaLIMS.