

PEAKS Q: Software for MS-based quantification of stable isotope labeled peptides

Weijie Yang¹, Weiwu Chen¹, Iain Rogers¹, Sean Bendall², Derek Smith², Bin Ma³, Gilles Lajoie³

¹Bioinformatics Solutions Inc., Waterloo, ON, ²Genome BC Proteomics Centre, Victoria, BC., ³University of Western Ontario, London, ON.



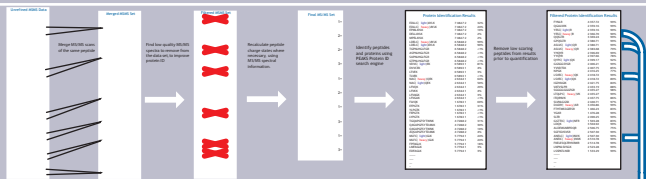
Introduction

Several mass spectrometry-based stable isotope labeling technologies have been developed for global proteome profiling. These include methods for *in vivo* labeling, such as ¹⁵N/¹³C and SILAC (Stable Isotope Labeling with Amino Acids in Cell Culture), and *in vitro* isotope labeling of target peptides at their N/C terminal or at specific residues. In this work we describe a new software, PEAKS-Q, designed to automatically identify and quantify proteins from these isotope labeling experiments. The software is written in Java and includes an intuitive graphical user interface.

Methods

The following functions are necessary:

- 1) Data Pre-processing: Sophisticated algorithms are used to merge spectra with similar retention time and m/z values, determine charge state when necessary, remove poor quality MS/MS scans, and remove noise, centroid and deconvolute data within MS/MS scans.
- 2) Protein identification: Peptides are identified from MS/MS data with the PEAKS protein identification algorithms.
- 3) Protein quantification: the ratio of each identified labeled peptide is calculated from the intensities of MS peaks that differ in mass by the mass of the label. The abundance of a peptide is obtained by averaging ratios from all observed charge state species of that peptide. Statistical methods are then applied to calculate the relative protein abundance and its associated standard deviation.



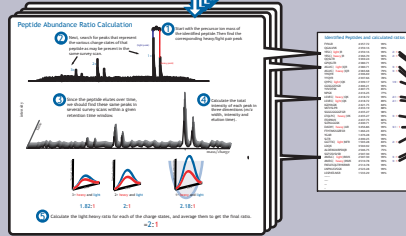
Removal of Deficient MS/MS spectra prior to analysis

Electrical noise, poor detection and contaminants scanned by the MS mean that only a small portion of data are quality MS/MS spectra representing peptides. Database search engines and *de novo* sequencing tools are adequate in discarding the bad spectra; nevertheless, false positives abound, and plenty of time is wasted. Hence a filter that eliminates poor spectra before the analysis can significantly improve throughput and robustness in a quantification software.

The algorithm included in the software had 99.61% accuracy in finding spectra of deficient quality and took only 130 seconds to reduce a

Peptide charge determination

Usually, we examine the initial MS survey scan of a peptide to determine the peptide survey scan charge state. But we cannot use this method with low resolution data, as most low-resolution instruments. If we let a protein identification tool decide the charge, we increase the risk of false positive matches, and triple the search time. As such, an algorithm is built into the software to find precursor charges with confidence, using low resolution tandem mass spectra data. The algorithm took less than four seconds to correctly assign charge on 315 spectra, with 92% accuracy².



Protein quantification

The ratio of each identified labeled peptide is calculated from the intensities of MS peaks that differ in mass by the mass of the heavy and light labels. The abundance of a peptide is obtained by averaging the ratios from all the observed charge state of that peptide. Dixon's test algorithm is used to remove extreme values (outliers) from a continuous data set. Statistical methods are then applied to calculate the relative protein abundance and its associated log deviation.

Let I denote peak intensity where I^L and I^H denote light and heavy peak intensities respectively. ϕ^Z denotes all isotopic peak groups where "Z" is the charge of the peptide. Assume a peptide is eluted between a range of retention

time range (T_1 to T_2) where $1 \leq k \leq n$, so that I_k^L and I_k^H can be written as follows:

$$I_k^L = \sum_{i=1}^n a_i \cdot 10^{i \cdot \log_{10}(1.00375)}$$

$$I_k^H = \sum_{i=1}^n a_i \cdot 10^{i \cdot \log_{10}(1.00375) + Z \cdot \log_{10}(1.00375)}$$

And the ratios (r) of peptide can be calculated by the following equation:

$$r = \frac{I_k^H}{I_k^L} = \frac{\sum_{i=1}^n a_i \cdot 10^{i \cdot \log_{10}(1.00375) + Z \cdot \log_{10}(1.00375)}}{\sum_{i=1}^n a_i \cdot 10^{i \cdot \log_{10}(1.00375)}}$$

Isotope Distribution

The software contains a new algorithm to predict the isotope distribution of a given peptide. The isotope abundance distribution is useful for comparing observed mass spectrometry data in helping to predict the number of atoms of a given element in the formula, and to distinguish signal peaks from chemical noise. The abundance distribution D is calculated using the observed natural abundance of each element in the formula, and convoluting these natural occurring abundances to predict the experimental isotope abundances.

Improved time complexity (for abundance distribution calculation):

$$D^Z(t) = \sum_{i=1}^n D^Z(t - i \cdot \log_{10}(1.00375)) \cdot 10^{i \cdot \log_{10}(1.00375)}$$

Experimental Results: BSA with cleavable ICAT

BSA samples were labeled with either light or heavy cleavable ICAT reagent and digested with trypsin. The light and heavy labeled samples were then mixed together and ratio of approximately 1:1.1, 2.4, 2.1 and 4:1. These samples were analyzed by LC-MS and MS/MS on a Waters QTRAP instrument. The software successfully identified 4 ICAT derived peptides that differ exactly by 9 Da as light/heavy pairs. The correct abundance ratio for each sample was determined, which indicates that the software can accurately determine abundance ratios over the dynamic range provided for the labeling experiment.

Abundance	Peptide	Source	Q Score	Protein	Ratio
1.00	ICAT-BSA(1)	BSA	1.00	BSA	1.00
0.78	GAGL-LPLK	BSA	0.78	BSA	0.78
0.76	VLCV-ANGOTLBSK	BSA	0.76	BSA	0.76
0.52	ICAT-BSA(2)	BSA	0.52	BSA	0.52
0.49	ICAT-BSA(3)	BSA	0.49	BSA	0.49
0.45	ICAT-BSA(4)	BSA	0.45	BSA	0.45
0.37	ICAT-BSA(5)	BSA	0.37	BSA	0.37
0.35	ICAT-BSA(6)	BSA	0.35	BSA	0.35
0.32	ICAT-BSA(7)	BSA	0.32	BSA	0.32
0.27	ICAT-BSA(8)	BSA	0.27	BSA	0.27
0.23	ICAT-BSA(9)	BSA	0.23	BSA	0.23
0.20	ICAT-BSA(10)	BSA	0.20	BSA	0.20
0.18	ICAT-BSA(11)	BSA	0.18	BSA	0.18
0.17	ICAT-BSA(12)	BSA	0.17	BSA	0.17
0.16	ICAT-BSA(13)	BSA	0.16	BSA	0.16
0.15	ICAT-BSA(14)	BSA	0.15	BSA	0.15
0.14	ICAT-BSA(15)	BSA	0.14	BSA	0.14
0.13	ICAT-BSA(16)	BSA	0.13	BSA	0.13
0.12	ICAT-BSA(17)	BSA	0.12	BSA	0.12
0.11	ICAT-BSA(18)	BSA	0.11	BSA	0.11
0.10	ICAT-BSA(19)	BSA	0.10	BSA	0.10
0.09	ICAT-BSA(20)	BSA	0.09	BSA	0.09
0.08	ICAT-BSA(21)	BSA	0.08	BSA	0.08
0.07	ICAT-BSA(22)	BSA	0.07	BSA	0.07
0.06	ICAT-BSA(23)	BSA	0.06	BSA	0.06
0.05	ICAT-BSA(24)	BSA	0.05	BSA	0.05
0.04	ICAT-BSA(25)	BSA	0.04	BSA	0.04
0.03	ICAT-BSA(26)	BSA	0.03	BSA	0.03
0.02	ICAT-BSA(27)	BSA	0.02	BSA	0.02
0.01	ICAT-BSA(28)	BSA	0.01	BSA	0.01

Experimental Results: Genome BC Proteomics Centre; E. Coli Complex mixture, ICAT

Proteins from *E. Coli* bacteria grown at the Genome BC Proteomics Centre were labeled with the isotopically light and heavy ICAT reagent and were analyzed by LC-MS/MS on an Applied Biosystems QSTAR instrument. The proteins were identified by PEAKS Studio Catalog database with the NCI95 database with the taxonomy defined as *E. coli*. Sixty five (65) proteins were identified from 1055 tandem mass spectra of ICAT labeled peptide pairs. Of all identified peptides, 98 peptides had acceptable chromatographic peaks in both the light and the heavy isotopic forms. The range of ratios determined by the software was between 0.73 and 0.86. About 80% of the peptide abundance ratios were within one log deviation. Sixteen (16) proteins had significant abundance changes, as measured by statistically derived t values. By evaluating the t values to specify the significance of protein abundance changes, the software is clearly capable of focusing and quantifying proteins of interest from a very large background.

Experimental Results: Genome BC Proteomics Centre; 5 to 1 BSA ICAT

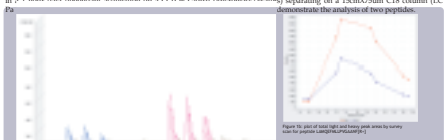
Two samples, both containing BSA, HELMANN, ALBU_BOVIN, and ICA, BSA were labeled - one with isotopically light ICAT reagent and the other with heavy ICAT - at the Genome BC Proteomics Centre. The two samples were combined together to create one mixture with 5:1 light to heavy ratio. The mixture was subsequently analyzed by LC-MS/MS on an ABI QSTAR (ESI-QTRAP), and the resulting data was processed for quantification. All three proteins were confidently identified by a PEAKS Protein ID search against the UniProt database. Table 1 summarizes the results of identification and

Abundance	Peptide	Source	Q Score	Protein	Ratio
1.00	ICAT-BSA(1)	BSA	1.00	BSA	1.00
0.78	GAGL-LPLK	BSA	0.78	BSA	0.78
0.76	VLCV-ANGOTLBSK	BSA	0.76	BSA	0.76
0.52	ICAT-BSA(2)	BSA	0.52	BSA	0.52
0.49	ICAT-BSA(3)	BSA	0.49	BSA	0.49
0.45	ICAT-BSA(4)	BSA	0.45	BSA	0.45
0.37	ICAT-BSA(5)	BSA	0.37	BSA	0.37
0.35	ICAT-BSA(6)	BSA	0.35	BSA	0.35
0.32	ICAT-BSA(7)	BSA	0.32	BSA	0.32
0.27	ICAT-BSA(8)	BSA	0.27	BSA	0.27
0.23	ICAT-BSA(9)	BSA	0.23	BSA	0.23
0.20	ICAT-BSA(10)	BSA	0.20	BSA	0.20
0.18	ICAT-BSA(11)	BSA	0.18	BSA	0.18
0.17	ICAT-BSA(12)	BSA	0.17	BSA	0.17
0.16	ICAT-BSA(13)	BSA	0.16	BSA	0.16
0.15	ICAT-BSA(14)	BSA	0.15	BSA	0.15
0.14	ICAT-BSA(15)	BSA	0.14	BSA	0.14
0.13	ICAT-BSA(16)	BSA	0.13	BSA	0.13
0.12	ICAT-BSA(17)	BSA	0.12	BSA	0.12
0.11	ICAT-BSA(18)	BSA	0.11	BSA	0.11
0.10	ICAT-BSA(19)	BSA	0.10	BSA	0.10
0.09	ICAT-BSA(20)	BSA	0.09	BSA	0.09
0.08	ICAT-BSA(21)	BSA	0.08	BSA	0.08
0.07	ICAT-BSA(22)	BSA	0.07	BSA	0.07
0.06	ICAT-BSA(23)	BSA	0.06	BSA	0.06
0.05	ICAT-BSA(24)	BSA	0.05	BSA	0.05
0.04	ICAT-BSA(25)	BSA	0.04	BSA	0.04
0.03	ICAT-BSA(26)	BSA	0.03	BSA	0.03
0.02	ICAT-BSA(27)	BSA	0.02	BSA	0.02
0.01	ICAT-BSA(28)	BSA	0.01	BSA	0.01

Note: 1) taxonomy of quantification results: HELMANN, ALBU_BOVIN, and ICA; 2) PEAKS Protein ID search details for GAGL-LPLK and VLCV-ANGOTLBSK.

Experimental Results: SILAC

HeLa cells were passaged 2 times (1:10 for one week in media containing 0.5% dialyzed FBS) followed by another 2x passage (DMEM containing either standard or heavy [¹³C]6 arginine (Specialty Media / Sigma). 90-95% labeling was confirmed after the 2nd passage by manual MS analysis. Cells were then harvested and mixed in 1:1, 1:4, and 4:1 ratios according to total cell count. Actual protein ratios may have varied from the intended ratios due to cell counting error ($\pm 1.35\%$) and varying levels of protein in each cell. A fourth, control sample, was prepared with [¹³C] Arg. Cells were lysed in RSM using 100mM ammonium bicarbonate, reduced with DTT, and alkylated with iodoacetamide prior to digest with trypsin (proteome). Digests were concentrated and desalted on a C18 solid phase extraction cartridge (Waters) and a small fraction of each digest was analyzed by LC-MS/MS on a 5.9 bore HPLC column (Agilent) with a 7500 G6 (Agilent) separating on a 150mm X 2.1mm C18 column (LC



Abundance	Peptide	Source	Q Score	Protein	Ratio
1.00	ICAT-BSA(1)	BSA	1.00	BSA	1.00
0.78	GAGL-LPLK	BSA	0.78	BSA	0.78
0.76	VLCV-ANGOTLBSK	BSA	0.76	BSA	0.76
0.52	ICAT-BSA(2)	BSA	0.52	BSA	0.52
0.49	ICAT-BSA(3)	BSA	0.49	BSA	0.49
0.45	ICAT-BSA(4)	BSA	0.45	BSA	0.45
0.37	ICAT-BSA(5)	BSA	0.37	BSA	0.37
0.35	ICAT-BSA(6)	BSA	0.35	BSA	0.35
0.32	ICAT-BSA(7)	BSA	0.32	BSA	0.32
0.27	ICAT-BSA(8)	BSA	0.27	BSA	0.27
0.23	ICAT-BSA(9)	BSA	0.23	BSA	0.23
0.20	ICAT-BSA(10)	BSA	0.20	BSA	0.20
0.18	ICAT-BSA(11)	BSA	0.18	BSA	0.18
0.17	ICAT-BSA(12)	BSA	0.17	BSA	0.17
0.16	ICAT-BSA(13)	BSA	0.16	BSA	0.16
0.15	ICAT-BSA(14)	BSA	0.15	BSA	0.15
0.14	ICAT-BSA(15)	BSA	0.14	BSA	0.14
0.13	ICAT-BSA(16)	BSA	0.13	BSA	0.13
0.12	ICAT-BSA(17)	BSA	0.12	BSA	0.12
0.11	ICAT-BSA(18)	BSA	0.11	BSA	0.11
0.10	ICAT-BSA(19)	BSA	0.10	BSA	0.10
0.09	ICAT-BSA(20)	BSA	0.09	BSA	0.09
0.08	ICAT-BSA(21)	BSA	0.08	BSA	0.08
0.07	ICAT-BSA(22)	BSA	0.07	BSA	0.07
0.06	ICAT-BSA(23)	BSA	0.06	BSA	0.06
0.05	ICAT-BSA(24)	BSA	0.05	BSA	0.05
0.04	ICAT-BSA(25)	BSA	0.04	BSA	0.04
0.03	ICAT-BSA(26)	BSA	0.03	BSA	0.03
0.02	ICAT-BSA(27)	BSA	0.02	BSA	0.02
0.01	ICAT-BSA(28)	BSA	0.01	BSA	0.01

If we can assume that the abundance ratio of any peptide belonging to a protein is representative of that protein's abundance ratio, then in theory it follows that all peptides belonging to a protein should have the same abundance ratio. An outlier is a relatively small or large data point within a data set where these values are statistically different from the main body of the data. Outliers should be removed when calculating the protein abundance.

To find and remove outliers, the software computes the ratio between the difference of the minimum [or maximum value] with its neighbor value and the difference of the maximum and minimum values. This ratio should follow a certain distribution. The outlying minimum [or maximum value] removed from the data set if it does not follow the supported distribution³.

of two which are less than n , then combine them together to reproduce n 's abundance.

The method is very fast, the result is very accurate, and the advantage is accuracy/efficiency when calculating distributions for extremely large proteins/peptides. For protein p[7460236], whose molecular mass is 1184139.08Da, and has 13288 amino acids, the calculation time was less than seven seconds.

Results and conclusions

PEAKS Q is demonstrated to be robust, easy to use software for quantification of heavy and light peptides. From mass spectrometry, its protein quantification accuracy is well within experimental error, and significance of abundance changes between samples is easy to see.

References

1. Chen, C., Rogers, J., Hillstedt, et al. MS/MS spectra of multiple peptides within database search. *Anal. Chem.* 2004, 76(12), 3000-3006.
2. Chen, C., Rogers, J., Moore, J. *Label-Free Protein Quantification from Top MS/MS Peaks*. *Anal. Chem.* 2005, 77(12), 3900-3906.
3. Smith, R., Kuster, B., Fernandez, M., Pandey, A., and Mann, M. (2003) *Biol Chem* 278, 3011-3016.
4. Pandey, A., Anderson, J., and Mann, M. (2003) *Science* 299, 171-176.
5. Smith, L., and Pandey, A. (2003) *Trends Biochem Sci* 28, 363-364.
6. Jensen, F. J. (2002) *Chemistry & Biology* 9, 107-110.
7. Hardman, M., and Righetti, P. G. (2003) *Mass Spectrom Rev* 23, 203-262.
8. Saito, S., Yan, Y. P., Whitmore, E. L., Garcia, L. A., Chen, H., Taitel, R. L., and Borchman, M. L. (1999) *Anal. Chem.* 71, 1000-1004.
9. Saito, S. A., and Whitmore, E. L. (2003) *Trends Cell Biol* 13, 67-73.
10. Hoopkin, J., and Whitmore, E. L. (2003) *Trends Cell Biol* 13, 67-73.
11. Perkins, D. N., Pappas, D. J., Coon, D. M., and Colubelli, J. S. (1999) *Electrophoresis* 20, 3033-3038.
12. Wang, Q., Man, P. J., and Hill, A. (1999) *Mol Cell Biol* 19, 3857-3860.
13. Holgado-Madruga, M., Enlis, D. R., Mowbray, D. G., Gadea, A. R., and Wong, A. K. (1996) *Nature* 379, 266-268.
14. White, M. F. (1998) *Mol Cell Biochem* 182, 3-11.
15. Cahill, M. H., Brodeur, T. G., and Robinson, D. J. (1995) *Cell Regul* 2, 969-974.
16. de Yathco, A. M., Berg, C. M., Lorenson, M. L., Mowbray, C. J., and Row, J. L. (1992) *Nature* 357, 405-408.
17. Anderson, J. (1998) *Mol Cell Biochem* 182, 34-40.
18. Houghton, J., O'Connell, A., Miyake, O., Xiao, H., Mann, J., and Aebersold, R. (2003) *Biochem J* 370, 1-14.
19. Mann, M., and Aebersold, R. (2004) *Mol Cell Proteomics* 3, 1173-1176.
20. Mann, M., Ong, S. C., Graczyk, M., Steen, H., Jensen, O. N., and Pandey, A. (2003) *Trends Biochem Sci* 28, 101-107.
21. Houghton, J., O'Connell, A., Miyake, O., Xiao, H., Mann, J., and Aebersold, R. (2003) *Biochem J* 370, 1-14.
22. Payne, D. B., Rounsaville, A. J., Martin, F., Erickson, R. H., and Shergill, P. W. (1999) *Nature* 397, 613-615.
23. Houghton, J., O'Connell, A., Miyake, O., Xiao, H., Mann, J., and Aebersold, R. (2003) *Biochem J* 370, 1-14.
24. Robinson, D. J., and Cahill, M. H. (1995) *Mol Cell Biochem* 129, 209-214.
25. David, B. W. (1993) *Statistical Treatment for Biologists: Quantitative Data: Critical Values of Student's *t* Test, Pearson and Chi-Square Tables and the *F* Test*. Cambridge, New York: Cambridge University Press.

The copyright of this article is the property of the copyright owner. All rights reserved. No part of this article may be reproduced without the prior permission of the copyright owner.