

Improving Protein Identification Sensitivity by Combining MS and MS/MS Information for Shotgun Proteomics Using LTQ-Orbitrap High Mass Accuracy Data

Bingwen Lu, Akira Motoyama, Cristian Ruse, John Venable, and John R. Yates III*

Department of Chemical Physiology, SR11, The Scripps Research Institute, La Jolla, California 92037

We investigated and compared three approaches for shotgun protein identification by combining MS and MS/MS information using LTQ-Orbitrap high mass accuracy data. In the first approach, we employed a unique mass identifier method where MS peaks matched to peptides predicted from proteins identified from an MS/MS database search are first subtracted before using the MS peaks as unique mass identifiers for protein identification. In the second method, we used an accurate mass and time tag method by building a potential mass and retention time database from previous MudPIT analyses. For the third method, we used a peptide mass fingerprinting-like approach in combination with a randomized database for protein identification. We show that we can improve protein identification sensitivity for low-abundance proteins by combining MS and MS/MS information. Furthermore, “one-hit wonders” from MS/MS database searching can be further substantiated by MS information and the approach improves the identification of low-abundance proteins. The advantages and disadvantages for the three approaches are then discussed.

Shotgun proteomics refers to the global analysis of the digested products of protein mixtures such as tissues, cells, or protein complexes.^{1,2} Multidimensional protein identification technology is a popular approach for shotgun proteomic as it combines high-resolution separation with tandem mass spectrometry (MS/MS).^{3,4} In general, protein mixtures are proteolytically reduced to peptides. Peptide separation by multidimensional high-pressure liquid chromatography is directly coupled to a tandem mass spectrometer followed by database searching using a computer algorithm such as SEQUEST,⁵ Mascot,⁶ or OMSSA.⁷ This approach has become a powerful method for identifying and quantifying proteins.^{8–11}

One inherent disadvantage to the shotgun protein identification method outlined above is the dependency on the acquisition of tandem mass spectra. For current technologies, it is still impossible to perform tandem mass spectrometry on every single ion in a chromatographic window with a ± 3 amu precursor isolation window, although the current generation of tandem mass spectrometers has increased acquisition of MS/MS by a factor of 5–10. This limitation leads to undersampling of complex peptide mixtures, and thus, usually one needs to perform multiple MudPIT analyses to increase acquisition of measurable peptide species.^{12,13} It would be desirable to utilize data from the first round of mass spectrometry (MS1 or simply MS) to supplement protein identification, and this usually requires high mass accuracy mass spectrometers.

One type of instrument that has high mass accuracy capability is the Fourier transform ion cyclotron resonance mass spectrometer (FT-ICR-MS). FT-ICR-MS has the ability to measure peptide masses at low ppm levels.¹⁴ This creates the possibility of using the mass of a single peptide as a unique identifier for protein identification provided that the mass of the amino acid composition of the peptide is unique in a database. By using this concept, Smith and colleagues have developed and refined an accurate mass tag (and later accurate mass and time tag, AMT) approach for MS-based high-throughput proteomics studies using FT-ICR-MS.^{14–17}

* To whom correspondence should be addressed. Tel: 858-784-8862. Fax: 858-784-8883. E-mail: jyates@scripps.edu.

- (1) McDonald, W. H.; Yates, J. R., III. *Curr. Opin. Mol. Ther.* **2003**, *5*, 302–309.
- (2) Aebersold, R.; Mann, M. *Nature* **2003**, *422*, 198–207.
- (3) Washburn, M. P.; Wolters, D.; and Yates, J. R., III. *Nat. Biotechnol.* **2001**, *19*, 242–247.
- (4) Peng, J.; Elias, J. E.; Thoreen, C. C.; Licklider, L. J.; Gygi, S. P. *J. Proteome Res.* **2003**, *2*, 43–50.

- (5) Eng, J.; McCormack, A.; Yates, J., III. *J. Am. Soc. Mass Spectrom.* **1994**, *5*, 976–989.
- (6) Perkins, D. N.; Pappin, D. J.; Creasy, D. M.; Cottrell, J. S. *Electrophoresis* **1999**, *20*, 3551–3567.
- (7) Geer, L. Y.; Markey, S. P.; Kowalak, J. A.; Wagner, L.; Xu, M.; Maynard, D. M.; Yang, X.; Shi, W.; Bryant, S. H. *J. Proteome Res.* **2004**, *3*, 958–964.
- (8) McDonald, W. H.; Yates, J. R., III. *Dis. Markers* **2002**, *18*, 99–105.
- (9) Peng, J.; Schwartz, D.; Elias, J. E.; Thoreen, C. C.; Cheng, D.; Marsischky, G.; Roelofs, J.; Finley, D.; Gygi, S. P. *Nat. Biotechnol.* **2003**, *21*, 921–926.
- (10) Wang, R.; Prince, J. T.; Marcotte, E. M. *Genome Res.* **2005**, *15*, 1118–1126.
- (11) Cantin, G. T.; Venable, J. D.; Cociorva, D.; Yates, J. R., III. *J. Proteome Res.* **2006**, *5*, 127–134.
- (12) Liu, H.; Sadygov, R. G.; Yates, J. R., III. *Anal. Chem.* **2004**, *76*, 4193–4201.
- (13) Durr, E.; Yu, J.; Krasinska, K. M.; Carver, L. A.; Yates, J. R.; Testa, J. E.; Oh, P.; Schnitzer, J. E. *Nat. Biotechnol.* **2004**, *22*, 985–992.
- (14) Zimmer, J. S.; Monroe, M. E.; Qian, W. J.; Smith, R. D. *Mass Spectrom. Rev.* **2006**, *25*, 450–482.
- (15) Lipton, M. S.; Pasa-Tolic, L.; Anderson, G. A.; Anderson, D. J.; Auberry, D. L.; Battista, J. R.; Daly, M. J.; Fredrickson, J.; Hixson, K. K.; Kostandarithes, H.; Masselon, C.; Markillie, L. M.; Moore, R. J.; Romine, M. F.; Shen, Y.; Strimatter, E.; Tolic, N.; Udseth, H. R.; Venkateswaran, A.; Wong, K. K.; Zhao, R.; Smith, R. D. *Proc. Natl. Acad. Sci. U.S.A.* **2002**, *99*, 11049–11054.

Strittmatter et al.¹⁸ also experimented with the use of LC Time-of-flight (TOF) MS for AMT tag protein identification, since current TOF instruments can provide resolutions larger than 10 000 and 2–5 ppm mass accuracies. For the AMT approach developed by Smith and colleagues using either FT-ICR-MS or TOF-MS, the protein identifications are obtained by matching the AMT tags to a database of peptides based on MS/MS information. Recently, an attempt was made to apply AMT to the more recent hybrid instruments such as the LTQ-FT or LTQ-Orbitrap.¹⁹

Another common MS-based method for protein identification is peptide mass fingerprinting (PMF).^{20–23} This method is most commonly applied to purified or highly enriched proteins often isolated by two-dimensional gel electrophoresis. PMF is most effective when attempting to identify single proteins, but an iterative approach has been used to improve the identification of simple mixtures of proteins.^{24,25} PMF has also been combined with two-dimensional chromatography of intact proteins in studies of cell lysates.²⁶ Giddings et al. developed an interesting variation of PMF to identify open reading frames in bacterial genomes.²⁷ Statistical or probability-based approaches for PMF analysis help to assess the quality of matches, but no PMF approaches have attempted to assess the false discovery rate associated with analysis.^{21,27–29}

The linear ion trap-orbitrap is a hybrid Fourier transform mass spectrometer that combines the efficiency and sensitivity of the linear ion trap with the high mass accuracy and high resolution of the orbitrap mass analyzer.^{30,31} The LTQ-Orbitrap has been shown to routinely achieve sub-5-ppm mass accuracy at a dynamic range of more than 5000³² and thus should be suitable for a protein identification approach using the accurate mass and time tags.

Here we investigated and compared three methods for shotgun protein identification by combining MS and MS/MS information

using LTQ-Orbitrap high mass accuracy data. The intent was to extend the amount of information that could be obtained from a MudPIT type of experiment and provide additional supporting information for MS/MS-based “one-hit wonders”. In the first approach, we employed a unique mass identifier method using m/z values that are not matched to peptides generated from proteins identified by MS/MS (Supporting Information, SI, Figure 1A). In the second approach, we used an accurate mass and time tag method by building a potential mass and time tag database from previous MudPIT analyses (SI Figure 1B). The third method used a PMF method using the m/z values that incorporated a randomized database to assess the false discovery rate (SI Figure 1C). We show an improvement in protein identification sensitivity of low-abundance proteins by combining MS and MS/MS information and add additional information to substantiate one-hit wonders.

Finally, the advantages and disadvantages for the three approaches are discussed.

EXPERIMENTAL PROCEDURES

Yeast Whole Cell Lysate Experiment. A protease-deficient *Saccharomyces cerevisiae* strain BJ5460³³ was purchased from American Type Culture Collection (Manassas, VA). The strain was grown to mid-log phase (OD 0.6) in YPD, and cells collected by centrifugation were lysed as described previously.³ The lysed cells were separated into three fractions (soluble and lightly and heavily washed), and the soluble fraction was used in this study. The soluble fraction of cells was digested by a method slightly modified from the one described previously.³ Urea was added to the soluble fraction of the cell lysate to denature the proteins. Proteins were then reduced with TCEP, alkylated using iodoacetamide, and subsequently digested with trypsin. The digestion process was stopped by adding formic acid to a final concentration of 1%. The protein digest was aliquoted and stored at $-80\text{ }^{\circ}\text{C}$ prior to use.

The protein digest was pressure-loaded onto a fused-silica capillary desalting column containing 5 cm of 5- μm Polaris C18-A material (Metachem, Ventura, CA) packed into a 250- μm -i.d. capillary with a 2- μm filtered union (UpChurch Scientific, Oak Harbor, WA). The desalting column was washed with buffer containing 95% water, 5% acetonitrile, and 0.1% formic acid. After desalting, a 100- μm -i.d. capillary with a 5- μm pulled tip packed with 10 cm of 3- μm Aqua C18 material (Phenomenex, Ventura, CA) followed by 3 cm of 5- μm Partisphere strong cation exchanger (Whatman, Clifton, NJ) was attached to the filter union and the entire split column (desalting column–filter union analytical column) was placed inline with an Agilent 1100 quaternary HPLC (Palo Alto, CA) and analyzed using a modified 13-step separation described previously.³ The buffer solutions used were 5% acetonitrile/0.1% formic acid (buffer A), 80% acetonitrile/0.1% formic acid (buffer B), and 500 mM ammonium acetate/5% acetonitrile/0.1% formic acid (buffer C). Step 1 consisted of a 100-min gradient from 0 to 100% buffer B. Steps 2–12 had the following profile: 3 min of 100% buffer A, 2 min of $X\%$ buffer C, a 10-min gradient from 0 to 15% buffer B, and a 97-min gradient from 15 to 45% buffer B. The 2-min buffer C percentages (X) were 10, 15, 20, 25, 30, 35, 40, 45, 50, 60, and 70%, respectively, for step 2 to step 12 analyses. For the final step, the gradient contained: 3 min of 100% buffer

- (16) Belov, M. E.; Anderson, G. A.; Wingerd, M. A.; Udseth, H. R.; Tang, K.; Prior, D. C.; Swanson, K. R.; Buschbach, M. A.; Strittmatter, E. F.; Moore, R. J.; Smith, R. D. *J. Am. Soc. Mass Spectrom.* **2004**, *15*, 212–232.
- (17) Pasa-Tolic, L.; Masselon, C.; Barry, R. C.; Shen, Y.; Smith, R. D. *Biol. Tech.* **2004**, *37*, 621–639.
- (18) Strittmatter, E. F.; Ferguson, P. L.; Tang, K.; Smith, R. D. *J. Am. Soc. Mass Spectrom.* **2003**, *14*, 980–991.
- (19) May, D.; Fitzgibbon, M.; Liu, Y.; Holzman, T.; Eng, J.; Kemp, C. J.; Whiteaker, J.; Paulovich, A.; McIntosh, M. J. *Proteome Res.* **2007**, *6*, 2685–2694.
- (20) Yates, J. R., III; Speicher, S.; Griffin, P. R.; Hunkapiller, T. *Anal. Biochem.* **1993**, *214*, 397–408.
- (21) Pappin, D. J.; Hojrup, P.; Bleasby, A. J. *Curr. Biol.* **1993**, *3*, 327–332.
- (22) James, P.; Quadroni, M.; Carafoli, E.; Gonnet, G. *Biochem. Biophys. Res. Commun.* **1993**, *195*, 58–64.
- (23) Henzel, W. J.; Billeci, T. M.; Stults, J. T.; Wong, S. C.; Grimley, C.; Watanabe, C. *Proc. Natl. Acad. Sci. U.S.A.* **1993**, *90*, 5011–5015.
- (24) Jensen, O. N.; Podtelejnikov, A. V.; Mann, M. *Anal. Chem.* **1997**, *69*, 4741–4750.
- (25) Lim, H.; Eng, J.; Yates, J. R., III; Tollaksen, S. L.; Giometti, C. S.; Holden, J. F.; Adams, M. W.; Reich, C. I.; Olsen, G. J.; Hays, L. G. *J. Am. Soc. Mass Spectrom.* **2003**, *14*, 957–970.
- (26) Lubman, D. M.; Kachman, M. T.; Wang, H.; Gong, S.; Yan, F.; Hamler, R. L.; O’Neil, K. A.; Zhu, K.; Buchanan, N. S.; Barder, T. J. *J. Chromatogr., B* **2002**, *782*, 183–196.
- (27) Giddings, M. C.; Shah, A. A.; Gesteland, R.; Moore, B. *Proc. Natl. Acad. Sci. U.S.A.* **2003**, *100*, 20–25.
- (28) Eriksson, J.; Fenyo, D. *Proteomics* **2002**, *2*, 262–270.
- (29) Henzel, W. J.; Watanabe, C.; Stults, J. T. *J. Am. Soc. Mass Spectrom.* **2003**, *14*, 931–942.
- (30) Makarov, A. *Anal. Chem.* **2000**, *72*, 1156–1162.
- (31) Schwartz, J. C.; Senko, M. W.; Syka, J. E. *J. Am. Soc. Mass Spectrom.* **2002**, *13*, 659–669.
- (32) Makarov, A.; Denisov, E.; Lange, O.; Horning, S. *J. Am. Soc. Mass Spectrom.* **2006**, *17*, 977–982.

(33) Jones, E. W. *Methods Enzymol.* **1991**, *194*, 428–453.

A, 20 min of 100% buffer C, a 10-min gradient from 0 to 15% buffer B, and a 107-min gradient from 15 to 70% buffer B.

As peptides eluted from the microcapillary column, they were electrosprayed directly into an LTQ-Orbitrap mass spectrometer (ThermoFisher, San Jose, CA) with the application of a distal 2.5-kV spray voltage. A cycle of one full FT scan mass spectrum (400–2000 m/z , resolution of 60 000) followed by five data-dependent MS/MS acquired in the linear ion trap with normalized collision energy (setting of 35%) was repeated continuously throughout each step of the multidimensional separation. Application of mass spectrometer scan functions and HPLC solvent gradients were controlled by the XCalibur data system.

Overall Data Analysis Outlines. The general outline is shown in SI Figure 1. We investigated three data analysis flows to more fully utilize the data generated by the new generation of hybrid mass spectrometers that produce large numbers of tandem mass spectra and high mass accuracy precursor ions. In the first approach (SI Figure 1A), we employed a unique mass identifier method where m/z values that match to proteins identified from the MS/MS database search are first subtracted before the m/z values are used as unique mass identifiers for protein identification. In the second method (SI Figure 1B), we used an accurate mass and time tag method by building a mass and retention time database from previous MudPIT analyses. For the third method (SI Figure 1C), we used a peptide mass fingerprinting approach together with randomized databases using the m/z values obtained from high mass accuracy precursor scans of a regular MudPIT analysis. The details of these three approaches are given below.

Approach One: MS Peak Subtraction and Unique Identifiers. MS and MS/MS data are collected on an LTQ-Orbitrap mass spectrometer in the format of the instrument's raw file. Each raw mass spectrometry data file obtained from the LTQ-Orbitrap is converted to an MS2 file using RawExtract, an in-laboratory software program (SI Figure 1A),³⁴ which contains the MS/MS measurements. Each raw file is also processed by using the software ICR-2LS (Anderson, G.A., <http://ncrr.pnl.gov/software/>) to give a list of deisotoped MS measurements (m/z 's). A regular database search using the SEQUEST algorithm is then performed on MS2 (MS/MS data) to obtain a list of identified proteins P_1 , as described in the section Analysis of Tandem Mass Spectra. The deisotoped m/z 's are then matched to predicted tryptic peptides (allowing up to 1 internal R or K, see SI Figure 2) from the list of identified proteins P_1 . The remaining m/z 's are further matched to predicted singly modified tryptic peptides from the identified proteins set, P_1 , and then subtracted from the m/z list. We considered common post-translational modifications (SI Table 1) including the following: (1) phosphorylation (serine, threonine, tyrosine); (2) N-terminal acetylation; (3) oxidation (methionine, tryptophan); (4) lysine acetylation; (5) arginine monomethylation; (6) arginine dimethylation. Finally, the remaining m/z 's are used as unique identifiers by matching to predicted tryptic peptides from the remaining proteins in the database. This will generate a second list of identified proteins P_2 . The two protein lists P_1 and P_2 are then pooled together to give the combined identified proteins P_T .

Approach Two: MRT Database Construction and Accurate Mass and Time Tagging. Tandem mass spectra data from previous replicate MudPIT runs were analyzed using the procedure outlined in the section Analysis of Tandem Mass Spectra, with the additional action of recording the retention time for each identified peptide. To be able to assess false positive rates for the AMT method, reverse peptide identifications are built into the potential MRT database. The protein identification false positive rate is set at 5% for each individual MudPIT experiment.

After the MRT database is constructed, accurate mass and time tagging is performed as follows. Each raw file is processed by using the software ICR-2LS to give a list of deisotoped MS measurements (m/z 's). The deisotoped MS peaks are then used as accurate mass and time tags for protein identification by matching to the potential MRT database. Under this experiment, we allowed a 5 ppm mass tolerance and 5% retention time tolerance when performing accurate mass and time tagging. AMT will generate a list of identified proteins P_2 . The MS/MS data are processed as usual to give the list of identified proteins P_1 . The two protein lists P_1 and P_2 are pooled together to give the combined identified proteins P_T .

Approach Three: Peptide Mass Fingerprinting Using MudPIT-Based MS Data. In this approach, we use a randomized protein database to analyze the MS data for protein identification by peptide mass fingerprinting under a specific false positive rate. For protein identification from tandem mass spectra, estimation of false positive rates using a reversed protein database is becoming a standard method.⁴ However, for protein identification using peptide mass mapping, randomized protein databases are not used to assess false discovery rates but may be a good choice as we show below.

We first describe the procedure for constructing the random databases R_1 and R_2 . We then give the rationales for using these two random databases.

Randomized databases are generated as follows (SI Figure 3A). The target protein database is first read, and the amino acid frequencies and protein length information are recorded. Two randomized protein sequence databases are then generated, both conserving the amino acid frequencies of the target protein database. The first randomized protein database R_1 is composed of N (N should be reasonably large; we choose $N = 10\,000$) randomized proteins, each of length L_0 ($L_0 = 500$ in this study). The second randomized database R_2 has the same number of protein entries as the target database, with each protein entry having the same length for its corresponding entry in the target database. The first randomized database R_1 is used to estimate μ and σ of the number of m/z hits of the collected m/z list to a randomized protein. R_2 is used to estimate false discovery rate (SI Figure 3B).

We show below that the random database R_1 can be used to estimate the significance of an m/z -peptide match. For each MudPIT run, all deisotoped MS1 peaks are first matched to predicted tryptic peptides from the randomized database R_1 , allowing up to one internal tryptic site for each predicted tryptic peptide (SI Figure 3B). For each protein entry in R_1 , the number of redundant matching MS1 peaks (peaks matched to the same tryptic peptides within the protein entry, referred to as "redundant MZ hits") and unique matching peaks (peaks matched to different

(34) McDonald, W. H.; Tabb, D. L.; Sadygov, R. G.; MacCoss, M. J.; Venable, J.; Graumann, J.; Johnson, J. R.; Cociorva, D.; Yates, J. R., III. *Rapid Commun. Mass Spectrom.* **2004**, *18*, 2162–2168.

tryptic peptides within the protein entry, referred as “unique MZ hits”) are recorded. The distribution of unique MZ hits for replicate 1 follows a normal distribution (SI Figure 4A) while the distribution of redundant MZ hits does not (SI Figure 4B). Because the distribution of unique MZ hits to R_1 is a normal distribution, the significance of the actual number of unique MZ hits can be assessed according to the mean and standard deviation of the distribution of unique MZ hits. This method utilizes the principle for the classical Z-test. For the classical Z-test, we assume the population is a normal distribution. If we see a sample data point, we can test whether this data point is from the population by computing the Z-score for the data point according to the mean and standard deviation of the population. Here the distribution of unique MZ hits to random database R_1 serves as the population distribution. The significance of the unique MZ hits for each protein can then be assessed according to the mean and standard deviation of the distribution of unique MZ hits to R_1 . For this purpose, the mean and standard deviation of the number of unique MZ hits are calculated. For replicate 1 with a 5 ppm mass tolerance, the mean of unique MZ hits is 42.8789 with a standard deviation of 8.4879.

After the mean and standard deviation of the number of m/z hits of the collected m/z list to a randomized protein were obtained, we used the following protein identification scoring formula

$$z = \frac{\frac{x}{L} - \mu}{\sigma \sqrt{\frac{L_0}{L}}} \quad (1)$$

where x is the number of nonredundant m/z hits, L is the length of the protein, L_0 is the common length of the randomized protein database R_1 , μ is mean of unique MZ hits to the randomized protein database R_1 , and σ is the standard deviation of unique MZ hits to the randomized protein database R_1 . The scoring performed a count of the number of unique MZ hits. The count was then normalized by the length of the protein. The normalized count was then compared to the unique MZ hits distribution to R_1 to obtain the statistical Z-score. The Z-score itself can be used for assessing the significance of a protein identification.

A z score is calculated for each protein entry according to the number of unique MZ hits according to eq 1. The distributions of z scores for R_2 and the target database for replicate 1 are shown in Figure 1. The distribution of z scores for R_2 follows a standard normal distribution with mean z scores of 0.0298 (≈ 0) and a standard deviation of 0.9926 (≈ 1). However, the distribution of z scores for the target database does not follow a standard normal distribution. The mean z score for the target database is -0.0439 with a standard deviation of 1.6309—we have more protein entries in the two tails. By using a certain z score cutoff, we obtained a list of proteins from the target database and the randomized database R_2 . The false discovery rate can then be estimated by the proportion of protein hits to R_2 in this list of identified proteins.

Analysis of Tandem Mass Spectra. Tandem mass spectra were analyzed using the following software analysis protocol. MS/MS were searched with the SEQUEST algorithm⁵ against a yeast protein sequence database concatenated to a decoy database in

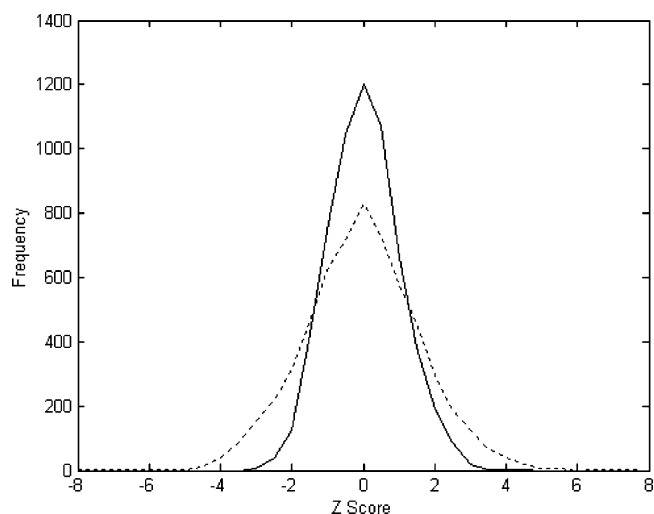


Figure 1. Scoring MZ hits using randomized databases (replicate 1). Blue solid line, distribution of z scores computed from unique MZ hits to the randomized yeast database R_2 . Red dash line, distribution of z scores computed from unique MZ hits to the yeast database (target database).

which the sequence for each entry in the original database was reversed.⁴ All searches were parallelized and performed on a Beowulf computer cluster consisting of 100 1.2 GHz Athlon CPUs.³⁵ No enzyme specificity was considered for any search. SEQUEST results were assembled and filtered using the DTA-Select (version 2.0) program.^{36,37} DTASelect 2.0 uses a quadratic discriminant analysis to dynamically set XCorr and ΔC_n thresholds for the entire data set to achieve a user-specified false positive rate. The false positive rates are estimated by the program from the number and quality of spectral matches to the decoy (reverse) database.

Retention Time Normalization. For multidimensional LC–MS/MS experiments, the retention time could be difficult to normalize. On the other hand, since we perform five MS/MS measurements following every MS measurement, peptide identifications from MS/MS experiments provide useful information for retention time normalization.

For experiment e_1 and experiment e_2 , each with a collection of peptide identifications from the SEQUEST search, the mean difference of retention time for all pairs of matching peptides between the two experiments can be calculated using the following formula:

$$\mu_t = \frac{1}{n} \sum_{i=1}^n t_{(p_{e1}, p_{e2})i} \quad (2)$$

where n is the number of matching peptide pairs between the two experiments, and $t_{(p_{e1}, p_{e2})i}$ indicates the retention time difference for the i th pair (p_{e1} , p_{e2}) of all the matching peptides. Using experiment e_1 as a reference, the retention times for experiment

(35) Sadygov, R. G.; Eng, J.; Durr, E.; Saraf, A.; McDonald, H.; MacCoss, M. J.; Yates, J. R., III. *J. Proteome Res.* **2002**, *2*, 211–215.

(36) Tabb, D. L.; McDonald, W. H.; Yates, J. R., III. *J. Proteome Res.* **2002**, *1*, 21–26.

(37) Cociorva, D.; Yates, J. R., III. DTASelect 2.0: Improving the Confidence of Peptide and Protein Identifications. In *54th ASMS Annual Meeting Proceedings*; 2006.

Table 1. SEQUEST ID Summary at 5% FP^a

	replicate 1	replicate 2	replicate 3
total proteins	1307	1372	1347
<i>P</i> ₁ proteins	447(76)	468(72)	453(72)
<i>P</i> ₂ proteins	860	904	894
FP	4.85%	4.96%	4.75%

^a. *P*₁ proteins: proteins identified by SEQUEST with one single peptide; *P*₂ proteins: proteins identified by SEQUEST with two or more peptides. FP: false positive rate. The number in the parentheses of *P*₁ proteins indicates the number of proteins that can be further verified by more peptides from the AMT approach (Approach Two) for each replicate experiment.

e_2 can then be normalized by adding μ_t to every retention time entry. The result of normalization of retention time of replicate 5 using replicate 4 as a reference is given in SI Figure 5.

EXPERIMENTAL RESULTS AND DISCUSSION

We performed six replicate 13-step MudPIT analyses on soluble yeast proteins. We used half of the data set (replicates 4, 5, and 6) for the construction of the potential MRT database and the other half (replicates 1, 2, and 3) for testing the scoring schemes. The MS and MS/MS information for replicates 1, 2, and 3 are given in SI Table 2. The false positive rates of identified proteins are controlled to be under 5%. Actual false positive rates might vary from case to case since sometimes it is impossible to control a false positive rate to be close to 5% under a specific filtering criterion.

Protein Identification by MS/MS Database Search. The first protein identification set *P*₁ is obtained from MS/MS database search by using SEQUEST followed by the validation program DTASelect 2.0. The protein identification false positive rate was set at 5% (at protein level) for each individual experiment for replicates 1, 2, and 3. Under the above filtering criteria, we identified 1307, 1372, and 1347 proteins for replicates 1, 2, and 3, respectively. When 2 peptides are required for an identification, we identified 860, 904, and 894 proteins from replicates 1, 2, and 3, respectively. When the data from replicates 1, 2, and 3 are pooled together and filtered by DTASelect2.0 requiring 2 peptides for each protein, we obtained 1119 proteins at a 0.44% false positive rate. The data are summarized in Table 1. For proteins identified by SEQUEST, roughly 34% of the proteins are identified by a single peptide (*P*₁, one-hit wonder).

Potential MRT Database Construction. Protein identifications by SEQUEST from replicates 4, 5, and 6 were used to construct the potential MRT database. To be able to assess the false positive rates of the accurate mass and time tagging approach, we built the reverse peptide/protein hits into the MRT database. The false positive rate was set at 5% (protein level) for each individual experiment. Table 2 gives a summary of the MRT database constructed from three replicate MudPIT runs. As estimated from reverse protein hits, the false positive rate for the whole potential MRT database was 1.38% at the peptide level and 7.29% at the protein level. The overall false positive rate (7.29%) was higher than each individual experiment (5%) since the true positive proteins tend to overlap while the false positive proteins tend to be different.

Protein Identification by Peak Subtraction and Unique Identifiers. From the set of proteins identified by SEQUEST, we

Table 2. MRT Database Summary

	no. forward	no. reverse	FP, %
copies	15516	130	0.84
peptides	8837	122	1.38
proteins	1673	122	7.29

generated theoretical tryptic peptides by requiring fully tryptic termini and allowing up to one internal tryptic site (see SI Figure 2). Theoretical (M + H)⁺s are then calculated for each theoretical peptide. From the set of deisotoped *m/z*'s obtained from MS1 measurements, we can calculate observed (M + H)⁺ by using the formula $MPlusH = mz \cdot z - z + 1$, where *mz* indicates the measured *m/z* and *z* indicates the charge of the ion. The observed (M + H)⁺s are then matched to theoretical (M + H)⁺s by allowing a given ppm mass tolerance. The matched (M + H)⁺s are removed from the set of observed (M + H)⁺s.

For the generation of predicted modified tryptic peptides, the nine common post-translational modifications in SI Table 1 are considered. The N-terminal acetylation refers to the N-terminal modification of proteins, not peptides. For each tryptic peptide, all nine modifications are considered, but the number of modifications is restricted to one modification for each predicted modified peptide. For example, if the tryptic peptide is AMLLTEDQK and the peptide is located at the beginning of the protein, then the following modifications will be considered and we will have four modified peptides (one for each modification): (1) N-terminal acetylation; (2) oxidation on M; (3) phosphorylation on T; (4) acetylation on K.

Theoretical (M + H)⁺s are then calculated for these predicted modified tryptic peptides and matched to the observed (M + H)⁺s by allowing a given ppm mass tolerance. The matched (M + H)⁺s are then removed from the observed (M + H)⁺s.

The remaining observed (M + H)⁺s are then used as unique identifiers for protein identifications by matching them to predicted tryptic peptides from the unidentified proteins in the database. To safeguard from false positive identifications, a selection criterion is applied to obtain the final list of protein identifications: a protein must be identified by at least by two unique *m/z* measurements (two distinct peptides) to be included in the protein identification set *P*₂. By applying these procedures using a 5 ppm mass tolerance, we identified 28, 63, and 28 proteins (*P*₂) from replicates 1, 2, and 3, respectively. The statistical details used for peak matching are given in SI Table 3.

Protein Identification by Accurate Mass and Time Tagging. We explored accurate mass and time tagging as outlined in SI Figure 1B. For this approach, all of the deisotoped MS1 peaks were used to search for matching peptides and proteins from the potential MRT database. A peak is disqualified for peptide and protein identification if it matches to two or more different peptides in the potential MRT database under a specified mass tolerance (5 ppm in this study) and retention time tolerance (5% in this study). More proteins and peptides are identified with increasing mass tolerance. However, the false positive rates also increase with relaxation of mass tolerance. We used a 5 ppm mass tolerance as the orbitrap mass analyzer can routinely achieve 5-ppm mass accuracy.³² We found that the false positive rates increased gradually when the retention time tolerance is relaxed, as shown

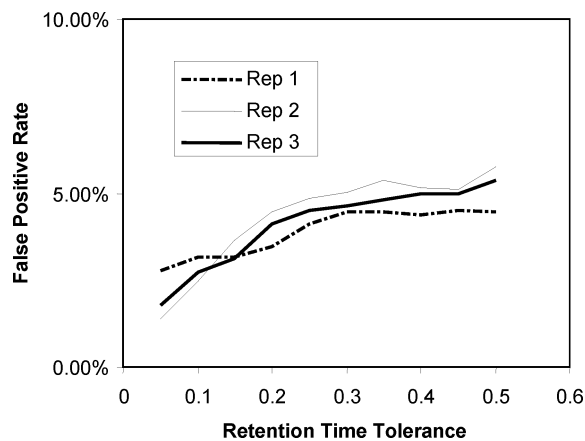


Figure 2. Effect of retention time tolerance on protein identification using the accurate mass and time tag approach.

in Figure 2. A 5% retention time tolerance provided a reasonable false positive rate.

At 5 ppm mass tolerance and 5% retention time tolerance, we obtained 922, 1007, and 973 protein identifications at the false positive rate of 3.15, 2.48, and 2.67%, for replicates 1, 2, and 3, respectively. When pooled together, these three replicates generated 1196 protein identifications containing 4.68% false positives.

When the proteins identified by AMT were combined with SEQUEST search results, we found that AMT results can further substantiate SEQUEST protein identifications. For SEQUEST-identified proteins with single peptides (one-hit wonders), 17.00 (76 out of 447), 15.38 (72 out of 468), and 15.89% (72 out of 453) of the proteins received additional peptide identifications from AMT, for replicates 1, 2, and 3, respectively (Table 1). When the AMT-identified proteins and SEQUEST-identified proteins are pooled together and then filtered with the requirement of at least two distinct peptide identifications for each protein, we obtained 1240 proteins with a false positive rate of 0.48% (i.e., 6 decoy matches). This is not achievable by using MS/MS data or MS data alone, but by combining MS and MS/MS information, we can improve protein identification sensitivity to a high confidence level.

Protein Identification by Using Peptide Mass Mapping.

We also investigated the use of randomized protein databases for the estimation of false positive rate for protein identification using MS mass mapping (SI Figure 1C). As described in the Experimental Procedures section, the target protein database is first read, and the amino acid frequencies and protein length information are recorded. Two randomized databases R_1 and R_2 are then generated (SI Figure 3A). Matches of MS1 peaks from each MudPIT runs to the target database as well as the randomized databases R_1 and R_2 are recorded.

A z score is then calculated for each protein entry according to the number of unique MZ hits according to eq 1. Using a z score cutoff of 3.0 (or chance probability of 0.001349), we identified 192 proteins from the target database and 9 proteins from the decoy database R_2 for replicate 1. Thus, the false positive rate at this z score cutoff for this specific experiment is 4.21%. Under the same z score cutoff of 3.0, we identified 194 and 170 proteins from replicates 2 and 3, respectively. By varying the z score cutoff, we can obtain protein identifications for different false positive rates, as shown in Table 3 for replicate 1 and SI Table 4 for replicates 2, and 3.

We found that the actual false positive rate is very close to the predicted false positive rate, as shown in Table 3. The column Predicted Random ID is calculated by the product of the chance probability and the number of protein entries in the randomized database R_2 (5996 as in this study). This further substantiates the validity of using randomized databases for the assessment of false positive rate using a PMF approach.

A validation of our PMF-like approach for protein identification using MudPIT data from a mixture of 17 known proteins was also supplied as Supporting Information to this paper. A six-step MudPIT generated 298 068 deisotoped MS peaks and 69 695 MS/MS spectra. We constructed a target database by appending the 17 proteins to the yeast protein database. MS/MS database search by SEQUEST followed by DTASelect identified 14 of the known proteins. By applying the PMF-like approach, we identified 7 proteins at the z score cutoff of 3.0. Among these 7 proteins, the top 6 proteins (as ranked by z scores) are from the 17 known proteins, verifying the reliability of the proteins identified by the scoring system.

The MudPIT-based PMF scoring method we showed here is relatively simple since we only consider the number of peptide matches and the length of the protein. The scoring might be further improved with more sophisticated schemes by using other information about the nature of each match, such as peptide length or missed tryptic sites.

Comparison of Different Approaches. Comparison of protein identification using different scoring schemes is given in Table 4. SEQUEST search of MS/MS data identified 1566 proteins at 4.98% protein level false positive rate from three replicate MudPIT runs replicates 1, 2, and 3. By using MS peaks not matched to proteins identified from MS/MS database search as unique identifiers (Unique Identifier, Approach One), we identified 92 additional proteins.

For peptide mass fingerprinting (PMF, Approach Three) using high mass accuracy MudPIT 2D-LC/MS data, we used the concept of decoy databases. By using randomized databases, we are able to establish a false positive rate for protein identification using PMF. By using this approach, we identified 271 proteins at 3.32% false positive rate by using PMF alone from the three MudPIT runs. Among these 271 proteins, 114 proteins are also identified by the SEQUEST search, with 157 new proteins identified by PMF.

The big overlap of identified proteins between PMF and SEQUEST provided additional statistical validation of the PMF method. Under a hypergeometric model, with a population size $N = 5996$ (the size of the number of proteins in the yeast database) and number of successes $k = 1566$ in population (number of proteins identified by SEQUEST), a sample size $n = 271$ (number of proteins identified by PMF) and number of successes $x = 114$ in sample (overlap between SEQUEST and PMF), according to the following hypergeometric formula

$$P(x) = \frac{\binom{k}{x} \binom{N-k}{n-x}}{\binom{N}{n}}$$

the probability of the null hypothesis that the match occurred by chance is 3.83×10^{-9} . The null hypothesis here is that PMF and SEQUEST identified the same subset of proteins by chance. The 157 proteins identified by PMF but not SEQUEST are of lower

Table 3. Protein Identified from Replicate 1 by Using PMF with Randomized Databases under Various z Score Cutoffs

z scores cutoff	chance probability	protein ID yeast DB	protein ID random DB R2	actual FP, %	predicted random ID	predicted FP, %
3.5	0.000233	112	2	1.79	1.40	1.25
3	0.001349	192	9	4.69	8.09	4.21
2.8	0.002555	246	19	7.72	15.32	6.23
2.6	0.004661	315	44	13.97	27.95	8.87
2.4	0.008197	391	80	20.46	49.15	12.57
2.2	0.013903	483	122	25.26	83.36	17.26
2	0.02275	582	183	31.44	136.41	23.44
1.8	0.03593	708	266	37.57	215.44	30.43
1.6	0.054799	867	382	44.06	328.57	37.90
1.4	0.080757	1048	536	51.15	484.22	46.20
1.2	0.11507	1257	725	57.68	689.96	54.89
1	0.158655	1486	967	65.07	951.30	64.02

Table 4. Comparison of Different Scoring Schemes^a

	Rep 1	Rep 2	Rep 3	total	extra
SEQUEST	1307 (4.85)	1372 (4.96)	1347 (4.75)	1566 (4.98)	0
approach one (5 ppm)	28 (na)	63 (na)	28 (na)	105 (na)	92
approach two (5 ppm, 5% RT)	922 (3.15)	1007 (2.48)	973 (2.67)	1196 (4.68)	103
approach three (5 ppm, z >= 3.5)	112 (1.79)	101 (1.98)	83 (1.20)	271 (3.32)	157

^a Numbers in parentheses indicate the false positive rates (%) for the identified proteins. na, not available; total, total identified proteins; extra, extra proteins identified compared to SEQUEST MS/MS search. For approach one, at least two unique peptide matches are required to identify a protein.

protein abundance compared to those identified by both PMF and SEQUEST (see below).

For accurate mass and time tagging (AMT, Approach Two), at 5 ppm mass tolerance and 5% retention time tolerance, we identified 1196 proteins at 4.68% false positive rate. When these proteins were compared with the 1566 proteins by SEQUEST, 103 proteins were found to be unique. When the proteins from SEQUEST and AMT are pooled together and filtered with an additional criterion of at least two distinct peptide identifications for a protein identification, we obtained 1240 proteins with a false positive rate of 0.48% (or 6 decoys). This is unachievable by using either MS/MS database search or MS-based approaches alone. By combining MS and MS/MS information, we can improve protein identification sensitivity while maintaining high confidence.

In total, the three approaches identified 342 extra proteins (SI Figure 6). This is a 21.8% improvement compared to MS/MS identification alone. These extra proteins are found to be consistently of lower abundance. The average abundance for proteins identified from SEQUEST is 2.63×10^4 copies/cell,³⁸ while the average abundances are 2.87×10^3 , 3.95×10^3 , and 2.63×10^3 copies/cell for proteins identified by unique mass identifiers, AMT, and PMF, respectively. Especially, for the 271 proteins identified by PMF, the average abundance of the 114 proteins that overlap with those of SEQUEST is 4.78×10^4 copies/cell, which is on average 18-fold more abundant than those identified by PMF only. The data show that our methods improve identification of low-abundance proteins.

The Unique Identifier approach (Approach One) will probably work better when the genome of the studied organism is relatively

small.¹⁵ When the genome is large, it is difficult to use accurately measured masses as unique identifiers for protein identification, since the larger the genome, the more peptides will have the same (or close) precursor masses. The specificity of the AMT approach (Approach Two), however, can be increased through the use of retention time as additional information for protein identification. The AMT approach can also leverage identifications from previous experiments. Single MS/MS matches, one-hit wonders from a specific experiment, can be further substantiated by accurately measured precursor masses of previously identified peptides as we showed here. Identifications by the AMT approach, however, are only limited to previous experiments and thus are not suitable for discovering unknown peptides. The PMF-like approach (Approach Three), instead, can identify proteins previously not identified. As we showed here, more than half of the proteins identified by the PMF-like approach were not identified by MS/MS database search. These proteins consistently have lower abundance levels compared to those identified by MS/MS database searching. The PMF-like approach, however, cannot identify as many proteins as the MS/MS database search, as shown by Table 4 as well as our supplementary results on the 17-protein mixture. Nevertheless, the PMF-like approach could be a protein identification method complementary to MS/MS database search.

CONCLUSIONS

In this study, we investigated and compared three approaches to extend identifications from shotgun experiments by combining MS and MS/MS information using LTQ-Orbitrap high mass accuracy data. In the first approach, MS peaks are first subtracted before being used as unique mass identifiers for protein identification. In the second approach, we explored the use of LTQ-Orbitrap

(38) Ghaemmaghami, S.; Huh, W. K.; Bower, K.; Howson, R. W.; Belle, A.; Dephoure, N.; O'Shea, E. K.; Weissman, J. S. *Nature* **2003**, *425*, 737–741.

data for AMT tagging. For the third method, we introduce the concept of decoy (randomized) databases for large-scale peptide mass fingerprinting. Comparing the different methods, the approach of combining MS/MS database search and AMT tagging is most promising from our data. In conclusion, by combining MS and MS/MS information, we can improve protein identification sensitivity with high confidence. One-hit wonders from MS/MS database search can be further verified by MS information, and the approach improves the identification of low-abundance proteins.

ACKNOWLEDGMENT

B.L. is supported by CFFT computational fellowship BALCH05X5. J.R.Y. acknowledges support from NIH

5R01MH067880-02, NIH P41 RR11823-10, and NIH R01 HL079442. The authors thank Drs. Daniel Cociorva, Meng-Qiu Dong, Daniel McClatchy, and Tao Xu for helpful reading of the manuscript. The authors also thank the reviewers for providing very detailed and helpful comments and feedbacks.

SUPPORTING INFORMATION AVAILABLE

Additional information as noted in text. This material is available free of charge via the Internet at <http://pubs.acs.org>.

Received for review August 9, 2007. Accepted January 4, 2008.

AC701697W