

## Trade-Off between High Sensitivity and Increased Potential for False Positive Peptide Sequence Matches Using a Two-Dimensional Linear Ion Trap for Tandem Mass Spectrometry-Based Proteomics

Hongwei Xie and Timothy J. Griffin\*

*Department of Biochemistry, Molecular Biology and Biophysics, University of Minnesota, Minnesota 55455*

Received December 20, 2005

**Abstract:** Two-dimensional linear ion trap mass spectrometers are rapidly becoming the new workhorse instruments for shotgun proteomic analysis of complex peptide mixtures. The objective of this study was to compare the potential for false positive peptide sequence matches between a two-dimensional ion trap instrument and a traditional, three-dimensional ion trap instrument. Through the comparative analysis of a complex protein sample, we found that in order to minimize false positive sequence matches, sequence match scoring criteria must be more stringent for data from the two-dimensional ion trap compared to the three-dimensional ion trap data. Given this increased potential for false positives, we also investigated two potential filtering strategies to reduce the false positive matches for data derived from the two-dimensional ion trap, including trypsin enzyme cleavage filtering, and the addition of peptide physicochemical information as a constraint, specifically peptide isoelectric point. The results described here provide a cautionary tale to researchers, demonstrating the need for careful analysis of MS/MS data from this new class of ion trap instruments, as well as the effectiveness of trypsin enzyme cleavage filtering and peptide pI information in maximizing high confidence protein identifications from this powerful proteomic instrumentation.

**Keywords:** proteomics • peptide identification • database searching • tandem mass spectrometry • false positive rate • linear ion trap • peptide isoelectric point

### Introduction

Ion trap mass spectrometers are relatively low in cost, robust, and have proven highly effective for tandem mass spectrometry (MS/MS) based shotgun proteomics.<sup>1–4</sup> Recently, a new generation ion trap mass spectrometer was introduced,<sup>5</sup> which has a quadrupole made of four hyperbolic cross-sectional rods giving it a linear two-dimensional (2D) design, and providing a number of significant functional enhancements, including 15 times higher ion capacity, 3 times faster scan rate, and

improved detection and trapping efficiency. The performance of a 2D and a 3D ion trap mass spectrometer (LTQ and LCQ instruments, respectively, marketed by Thermo Finnigan) in shotgun proteomics was recently systematically compared.<sup>6</sup> A 4–6-fold increase in the number of peptides and proteins identified on the 2D ion trap mass spectrometer compared to the 3D instrument was observed, demonstrating the improved capabilities of the 2D ion trap.

One challenge in shotgun proteomics is reducing false positive sequence matches, which are incorrect sequence matches assigned scores by sequence database search programs such as Sequest<sup>7</sup> or Mascot<sup>8</sup> that exceed a chosen threshold score.<sup>9,10</sup> To address this problem, sophisticated statistical tools have been developed,<sup>11–13</sup> including a probabilistic scoring algorithm called Peptide Prophet<sup>14</sup> which statistically models each specific proteomic dataset and takes into account multiple criteria to provide a Probability score (*P* score) to each match. However, even with these more advanced scoring algorithms, false positive matches still remain a difficulty in interpreting results from shotgun proteomic studies.

The objective of this study was to compare the potential for false positive peptide sequence matches from MS/MS data in shotgun proteomics using a 2D or 3D ion trap mass spectrometer. To this end, we systematically compared the potential for false positive sequence matches from a LTQ 2D ion trap and a LCQ Classic 3D ion trap, analyzing identical peptide samples derived from the soluble protein fraction of whole human saliva and fractionated by free flow electrophoresis (FFE).<sup>15</sup> False positive rates of peptide sequence matches obtained from both instruments and determined by the program Sequest<sup>7</sup> were estimated using reverse database searching.<sup>9</sup> Two strategies were evaluated to reduce the higher potential of false positive matches observed from the 2D LTQ ion trap data, including trypsin cleavage filtering and the addition of the physicochemical constraint of peptide isoelectric point (pI). The results described here provide a warning to researchers of the need for careful analysis of MS/MS data obtained from 2D ion trap instruments, and some effective strategies to increase the confidence of data obtained in proteomic studies utilizing this instrumentation.

### Experimental Section

**Preparation of Saliva Proteins.** The procedures for saliva sample collection, storage and isolation of soluble saliva proteins were described in a previous paper.<sup>16</sup> Briefly, whole

\* To whom correspondence should be addressed. Department of Biochemistry, Molecular Biology and Biophysics, University of Minnesota, 6-155 Jackson Hall, 321 Church St. SE, Minneapolis, MN 55455. Tel: (612) 624-5249. Fax: (612) 624-0432. E-mail: tgriffin@umn.edu.

un-stimulated saliva was collected from a healthy female subject in the University of Minnesota Oral Medicine Clinic using a previously described protocol.<sup>17</sup> The collected whole saliva was immediately placed on ice and frozen at  $-80^{\circ}\text{C}$ . The frozen saliva sample was thawed on ice, and 1 mL was removed and centrifuged at  $25\,000 \times g$  and  $4^{\circ}\text{C}$  for 30 min. The supernatant was collected and quantified by using the BCA protein assay (Pierce), giving 1.05 mg of total soluble proteins.

**Preparation of Peptides.** The saliva supernatant was thawed on ice and immediately brought to 100 mM with HEPES, pH 8.0 and 5 mM with TCEP and incubated overnight with  $20\,\mu\text{g}$  of trypsin (Promega, Madison, WI) at  $37^{\circ}\text{C}$ . The resulting peptides were first concentrated and desalted using a reverse-phase Sep-Pak cartridge (Waters, Milford, MA) and dried by vacuum centrifugation. Then the peptides were separated into 96 fractions by FFE based on preparative isoelectric focusing (IEF) of the peptide mixture and collected into 96 microtiter plate wells as described.<sup>16,18</sup> Immediately after FFE separation, the pH of each FFE fraction was measured using a microelectrode (Accumet Combination Micro Electrode, Fisher Scientific). Two equal volume aliquots of peptides from each FFE fraction were prepared as previously described,<sup>16,18</sup> one for microcapillary liquid chromatography ( $\mu\text{LC}$ ) MS/MS analysis using the 3D LCQ Classic ion trap, the other using 2D LTQ ion trap (both marketed by Thermo Corporation, San Jose, CA).

**$\mu\text{LC}$ -ESI MS/MS Analysis.** All FFE fractions were further separated by online  $\mu\text{LC}$  using the exact same chromatography conditions for both instruments. For the LCQ,  $\mu\text{LC}$  separation was done using an Agilent 1100 binary HPLC system, coupled to the LCQ ion trap mass spectrometer. For the LTQ, all online  $\mu\text{LC}$  separations were done on an automated Paradigm MS4 system (Michrom Bioresources, Inc., Auburn, CA). For both instruments, the inline analytical capillary column ( $75\,\mu\text{m}$  i.d.  $\times$  12 cm) was home-packed using C18 resin ( $5\,\mu\text{m}$ , 200 Å Magic C18AG, Michrom BioResource, Auburn, CA) and Picofrit capillary tubing ( $75\,\mu\text{m}$  i.d.  $\times$  12 cm, New Objective, Cambridge, MA). Each sample was first concentrated and de-salted by loading in buffer A (0.1% formic acid in solution of 5% acetonitrile and 95% water) on a C18 precolumn. Peptides were eluted using a linear gradient of 10–35% buffer B (0.1% formic acid in solution of 95% acetonitrile and 5% water) over 60 min, followed by isocratic elution at 80% buffer B for 5 min to wash the column with a flow rate of  $0.25\,\mu\text{L}/\text{min}$  across the column.

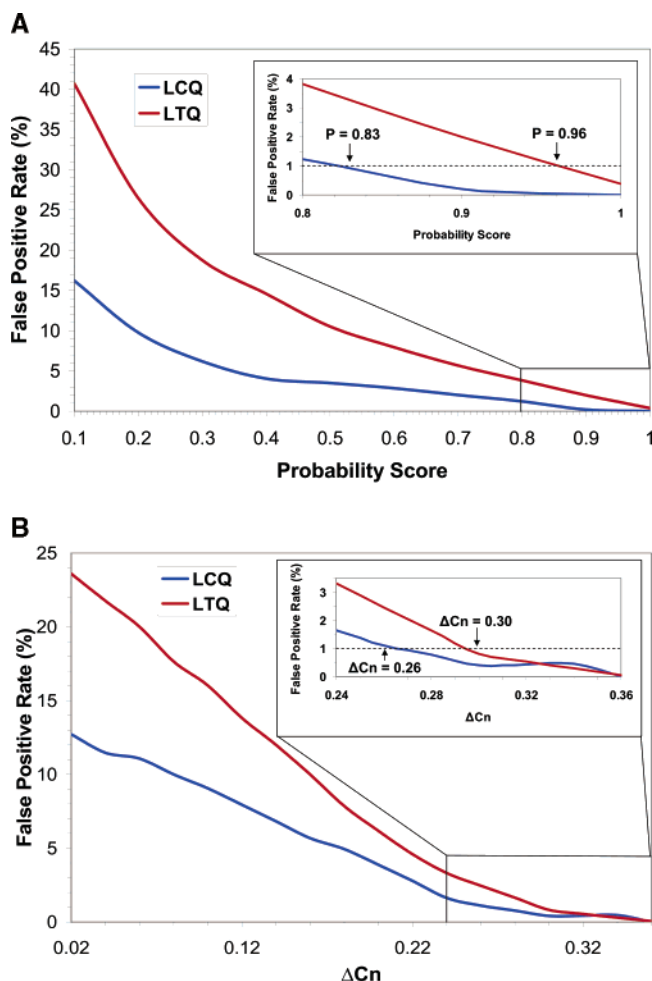
Peptides were analyzed by MS/MS using the LCQ or LTQ ion trap instruments using established optimal instrument parameter settings for proteomic experiments.<sup>6,19</sup> For the LCQ, peptides eluting from the capillary column were automatically selected for CID by the mass spectrometer using a protocol that alternated between one MS scan (summing  $3\,\mu\text{scans}$ ) and three MS/MS scans (each scan summing  $5\,\mu\text{scans}$ ) for the three most abundant precursor ions detected in the MS survey scan. Precursor  $m/z$  values selected for CID using a collision energy setting of 33% were dynamically excluded for 90 s after selection. The electrospray voltage was set to 1.7 kV. For the LTQ, ionized peptides eluting from the capillary column were selected for CID using a collision energy setting of 29% and a data-dependent procedure that alternated between one MS scan followed by four MS/MS scans for the four most abundant precursor ions in the MS survey scan. Both the MS and MS/MS spectra were acquired using a single  $\mu\text{scan}$  with a maximum fill-time of 50 milliseconds in the ion trap.  $m/z$  values selected for MS/MS were dynamically excluded for 30 s. The electrospray voltage was set to 2.0 kV. The operation of the both mass

spectrometers was controlled by the software Xcalibur. The mass range for precursor ion detection was set from 400 to 1800 Dalton on both instruments.

**Sequence Database Searching and Peptide Sequence Match Filtering.** The MS/MS spectra were sequence database searched using SEQUEST<sup>7</sup> (Thermo Finnigan, San Jose, CA) across a distributed cluster of 14 processors maintained at the University of Minnesota Supercomputing Institute. The MS/MS spectra were searched against a nonredundant human proteome sequence database from European Bioinformatics Institute (<http://www.ebi.ac.uk/IPI/IPIhuman.html>), with a reverse version of the same database attached at the end of the forward version. The search results were validated using Sequest scores and the publicly available peptide validation program Peptide Prophet ([www.systembiology.org/Default.aspx?pagename=proteomicssoftware](http://www.systembiology.org/Default.aspx?pagename=proteomicssoftware)),<sup>14</sup> which assigns a comprehensive probability score from 0 to 1 to each peptide identification based on its SEQUEST scores (Xcorr,  $\Delta\text{Cn}$ , Sp, RSp) and additional information of the peptide, including the mass difference between the precursor ion and assigned peptide, and the number of tryptic termini. The peptide sequence match results were organized and viewed using the software tool Interact.<sup>20</sup> The MS/MS spectra were first searched against the database with the enzyme trypsin specified, allowing matches to fully tryptic peptide sequences with up to two internal missed cleavage sites within the peptide sequence match; the data was also searched with no enzyme specified, and subsequently filtered according to expected tryptic cleavage sites, as described in the text. False positive rates were estimated using the algorithm originally described by Gygi and colleagues.<sup>9</sup> The predicted pI of peptide sequences was calculated according to Shimura and colleagues<sup>21</sup> using an automated script, and peptide pI values were automatically inputted into the Interact results file.

## Results and Discussion

The samples used in this comparative study were prepared from proteins contained in the soluble supernatant from a single whole human saliva sample. We have described the proteomic analysis of this saliva sample elsewhere.<sup>16</sup> For this present study, this sample represented an ideal, complex protein mixture obtained from a human source to compare the two ion trap instruments. Several groups have visualized the human saliva proteome via two-dimensional gels<sup>22–25</sup> and combined with our recent results<sup>16</sup> and others<sup>26</sup> it is known to be a bodily fluid with a complex proteome. Because it is from a human source, it also is a challenging sample for studying the phenomenon of false positive sequence matches, given the demonstrated potential for larger sequence databases (e.g., human) to produce increased false positive matches.<sup>10</sup> Saliva also provides easily collected, large amounts of total protein. As such, a single saliva sample was processed for this study, providing enough excess protein such that equal amounts of peptides from each FFE fraction could be analyzed separately on the 2D and 3D ion trap instruments. The saliva sample was handled on ice during collection and processing, and trypsin was immediately added to the isolated whole saliva supernatant for digestion, so as to minimize possible degradation of proteins by proteases in the saliva. The resulting peptide mixture was analyzed using a high-throughput shotgun proteomic strategy employing FFE as the first separation dimension, followed by  $\mu\text{LC}$ -ESI MS/MS.<sup>18</sup> For comparison, MS/MS spectra from equal aliquots of peptide mixtures from each FFE fraction were



**Figure 1.** Comparison of the false positive peptide sequence match rates between the LTQ and LCQ using different scoring methods for peptide sequence matches. A. *P* score assignments; B. Sequest score assignments, varying the  $\Delta Cn$  score, with the following constant Xcorr threshold values:  $\geq 2.0$  (+1 ions);  $\geq 2.5$  (+2 ions); and  $\geq 3.8$  (+3 ions). The tandem mass spectra were searched against the sequence database with trypsin specified cleavage.

acquired on the LTQ and the LCQ instruments, as described in the Experimental section. In total, 558 508 MS/MS spectra were obtained using the LTQ and 148 469 using the LCQ from the comparative analysis of peptides from total 64 FFE fractions prepared from the whole human saliva sample. A previously described<sup>9</sup> in-silico validation strategy for peptide sequence matches employing the chimeric database was used to estimate the false positive rate for any given scoring criteria of database searching results.

**LTQ Shows an Increased Potential for False Positive Matches versus the LCQ.** We first investigated differences in false positive rates between the two instruments using Peptide Prophet *P* score values for filtering the peptide sequence matches from the database search. For this initial search, trypsin cleavage was specified, constraining the possible matches to only tryptic peptides but allowing up to two internal missed cleavages. The plots of false positive rates versus *P* score values for both instruments are shown in Figure 1A. As expected, for both instruments the false positive rate gradually decreases with increasing *P* score. However, for any given *P* score, the false positive rate from the LTQ was significantly higher ( $\geq 2$ -fold)

than that from the LCQ, although the difference was decreased slightly with increasing *P* score value. To achieve a false positive rate below 1%, the *P* score value was 0.83 for the LCQ, and 0.96 for the LTQ. At an estimated false positive rate of 1%, 364 unique proteins could be identified from the tandem mass spectra acquired from the LTQ, but only 199 from the LCQ with the same confidence. Furthermore, nearly four times more total peptide sequences were matched from the LTQ, with two times more average peptide sequences matched per identified protein, compared to the LCQ, which is similar to the results from a previous report comparing the performance of these two instruments.<sup>6</sup>

Although the Peptide Prophet *P* score, which takes into account many different parameters in determining correct peptide matches,<sup>14</sup> helps to simplify the evaluation of peptide sequence matches, standard Sequest scores (Xcorr and  $\Delta Cn$ ) are still commonly used by many groups to evaluate peptide sequence matches and protein identifications. Therefore, we also compared false positive rates of peptide sequence matches derived from tandem mass spectra obtained from both the LTQ and LCQ at different values of the critical Sequest  $\Delta Cn$  score parameter. Since the optimal  $\Delta Cn$  score is affected by the chosen Xcorr threshold values, to simplify this comparison, we set the Xcorr threshold values at 2.0, 2.5, 3.8 for charge states of +1, +2, +3, respectively, which are stringent values which have been used by others.<sup>12,19</sup> We then calculated false positive rates at varying  $\Delta Cn$  score thresholds. As shown in Figure 1B, the LTQ again showed significantly higher false positive rates compared to the LCQ, especially at lower  $\Delta Cn$  values. To keep the false positive rate at 1% or below, a  $\Delta Cn$  threshold value of 0.26 was necessary for the LCQ, and 0.30 for the LTQ. These values are significantly higher than previously reported minimal values for  $\Delta Cn$ ,<sup>12,19</sup> for reasons that are described below, but still provide a means to compare the false positive rate between the two instruments. At an estimated false positive rate of 1%, 362 unique proteins could be derived from the tandem mass spectra acquired from the LTQ, and only 198 from the LCQ. Notably, filtering using  $\Delta Cn$  score coupled with Xcorr thresholds compared to filtering using *P* score for both the LTQ and LCQ instruments gave a similar number of identified proteins, peptides and average peptide sequence matches per protein, indicating these methods for scoring peptide matches are comparable in the identification of proteins.

**Use of Enzyme Cleavage and Physicochemical Peptide Information Maximizes High Confidence Protein Identifications from the LTQ.** The results shown in Figure 1 indicate a tradeoff to the increased sensitivity of the 2D ion trap in shotgun proteomics, which is an increased potential for false positive sequence matches from MS/MS data derived from this instrument. The main contributing factor to this increased false positive match rate is the increased number of MS/MS spectra acquired by the faster scanning 2D ion trap, which provides the increased sensitivity of this instrument in shotgun proteomic analyses, but also increases the potential for poor quality spectra, and spectra with unexpected fragmentation patterns, to erroneously match to peptide sequences and result in false positive matches. Given the increased potential for false positive sequence matches on the LTQ, we sought to investigate this phenomenon further, and to provide some potential solutions to aid researchers in maximizing the high confidence data that can be obtained from this instrumentation.

**Filtering by Protein Sequence Coverage.** The results shown in Figure 1 included single hits as well as proteins identified



from two or more unique peptide sequences. In shotgun proteomics, protein identifications derived from two or more unique peptide sequence matches are generally considered to be more confident than those derived from only a single peptide hit. Therefore, in some studies,<sup>19</sup> only those proteins derived from two or more unique peptide hits are accepted as correct. Although this approach provides a simple filtering method to ensure higher confidence for protein identifications, it sacrifices many correct protein identifications derived from high quality single hits. To compare the confidence of proteins identified from two or more unique peptide sequence matches from MS/MS spectra of these two different ion trap instruments in our present study, we first filtered the dataset used to generate Figure 1 above using a low stringency *P* score threshold of 0.2. We then only accepted protein identifications having two or more unique peptide matches. Using this criterion, the false positive rate for the LTQ was 2.61% (with 294 total proteins identified, reduced from 364 proteins identified when keeping single hits), whereas for the LCQ, the false positive rate was only 0.7% (with 105 total proteins identified, decreased from 199 proteins identified when keeping single hits). This demonstrates that for the LTQ there is still potential for false positive identifications, albeit it significantly decreased, even when single hit protein identifications are discarded. However, one also sacrifices many protein identifications when only considering proteins identified from multiple peptide matches to be correct. Interestingly, although the number of total proteins decreased, four new proteins were identified that were not identified using the initial strategy described above keeping all matches at a *P* score of 0.96 or above, indicating that filtering based on multiple peptide sequence matches can complement the use of threshold scoring values.

#### Filtering Using Trypsin Enzyme Cleavage Information.

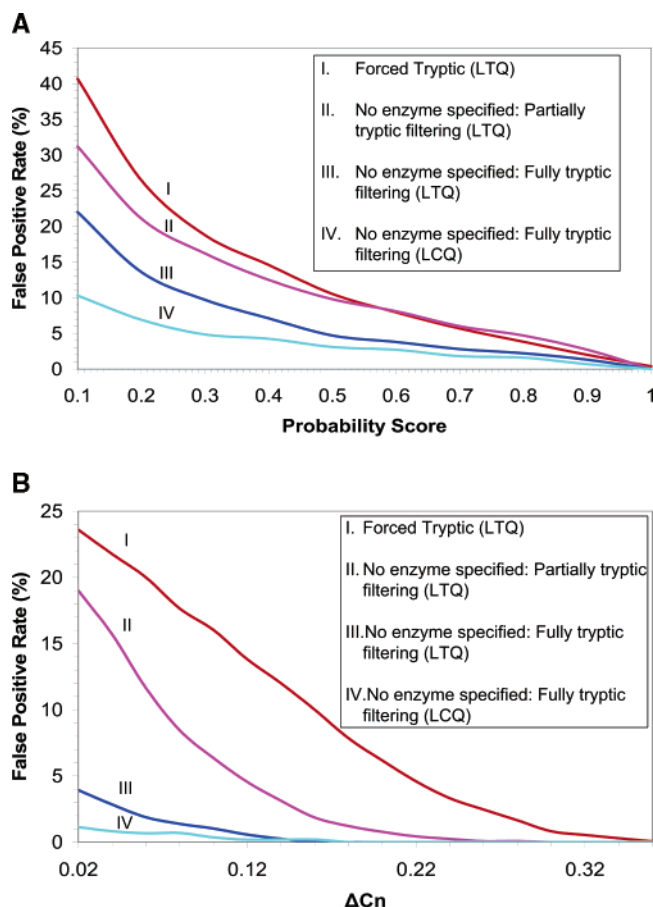
Given the limitation of sacrificing protein identifications from single-hit peptide sequence matches when considering only protein identifications from multiple peptide sequence matches to be correct, we next investigated the use of other strategies to maximize the number of high confidence protein identifications from 2D ion trap data. We first examined the use of trypsin enzyme cleavage information for filtering the peptide sequence matches derived from the LTQ to possibly decrease the potential for false positive peptide sequence matches. The results shown above in Figure 1 were derived using Sequest and searching with trypsin specified as the enzyme used in the database search parameters (i.e., "forced" tryptic peptide matches). This means that the search was conducted specifying that the peptide sequence matched to any MS/MS spectrum must be derived from proteolysis at either lysine or arginine, providing tolerance for up to two missed cleavage sites within the internal amino acid sequence of the peptide. This trypsin specification effectively constrains the number of possible peptide sequences matched to the experimental MS/MS data, and forces all sequence matches, whether correct or incorrect, to match to a tryptic peptide sequence. An alternate way in which to conduct the database search is to specify no enzyme constraint in the initial search parameters, allowing the program to match MS/MS spectra to any peptide sequence, regardless of enzyme cleavage sites in the protein. Post-sequence database search, the data is filtered based on the expectation that the peptides contained in the mixture should have been cleaved at specific amino acid sites (e.g., lysine and arginine for trypsin). This filtering strategy, based on what has been called the "trypsin distraction" effect,<sup>27</sup> effectively provides

a constraint on the data, as the probability that any given MS/MS spectra by chance matches to a tryptic peptide sequence versus all the other nontryptic possibilities is relatively low, providing a more confident assessment that matches to tryptic peptide sequences are indeed correct. This strategy effectively enables the filtering out of incorrect sequence matches, which will most likely match to nontryptic peptide sequences. This trypsin filtering approach has generally been taken in studies attempting to determine optimal Sequest scores (*Xcorr* and  $\Delta Cn$ ) for minimizing false positive rates.<sup>9,11,19</sup> As such, because our database search was first conducted using the forced trypsin constraint for the results shown in Figure 1, it is a main factor in our optimal  $\Delta Cn$  values for both instruments being quite a bit higher than these previous studies in which trypsin cleavage was not specified in the search. This increase in Sequest score thresholds has also been observed by others when using trypsin constrained searches.<sup>27</sup> Our results above also did not discard single hit peptide matches, different from these previous studies which specially treated these single hits in determining their optimal Sequest scoring values.<sup>12,19</sup>

To investigate the utility of enzyme constraints in reducing false positive matches, we re-searched the MS/MS data from both instruments with no enzyme specified, and then used expected tryptic cleavage as a subsequent filtering method for the sequence matches. After the database search, we then filtered the data in two different ways. First, we filtered allowing for peptide sequence matches derived from partially tryptic peptides, which are defined as peptides having at least one end (n-terminus or c-terminus) which has been cleaved by trypsin at an expected site.<sup>9,19</sup> These peptides may also have internal missed trypsin cleavage sights within the amino acid sequence of the peptide. Second, we filtered the data allowing for only fully tryptic peptides (i.e., peptide sequences with tryptic cleavage sites at both the n- and c-terminus) to be considered as correct matches. This effectively provides one more level of stringency, as the probability of poor MS/MS data matching to a fully tryptic peptide sequence by chance is even further decreased compared to partially tryptic peptide sequences.

Figure 2 shows the results using the enzyme filtering strategy, where the false positive rate is plotted against the assigned peptide match *P* values in Figure 2A; Figure 2B shows the same results, plotting the false positive rate versus assigned  $\Delta Cn$  values, using the same stringent values for *Xcorr* as described above. For comparison, curve I in both Figure 2A and B shows the results derived using the forced tryptic parameter in the database search, as shown in Figure 1. The number of proteins identified and total peptide sequence matches at an estimated false positive rate of 1% for each of the filtering methods shown in Figure 2 are summarized in Table 1.

There are several interesting observations from the results when using the trypsin enzyme constraint in filtering the data. For the case of using *P* score values for determining correct sequence matches (Figure 2A), filtering by allowing partially tryptic peptides (curve II) does not significantly decrease the false positive rate compared to the results when specifying trypsin in the search parameters (curve I). However, it does increase the total number of proteins identified, adding 49 proteins (Table 1). These additional protein identifications are derived from peptide sequences that are partially tryptic, and were not identified in our initial search when constraining the results to sequences with tryptic termini when specifying trypsin in the database search parameters. Given the sample being analyzed, this is not an entirely unexpected result, as it



**Figure 2.** False positive peptide sequence match rates for the LTQ using different database search and filtering methods. **A.**  $P$  score assignment; **B.** Sequest score assignment, varying the  $\Delta Cn$  score, with the following constant Xcorr threshold values:  $\geq 2.0$  (+1 ions);  $\geq 2.5$  (+2 ions); and  $\geq 3.8$  (+3 ions). The curves in each figure show results using different database search parameters and filtering, as described in the text.

**Table 1.** Summary of LTQ Results Using Different Database Searching and Filtering Parameters

scoring method <sup>a</sup>	filtering method	proteins	peptides
$P$ Score	forced tryptic	364	3893
	partially tryptic <sup>b</sup>	413	4882
	fully tryptic <sup>b</sup>	372	3121
$\Delta Cn$	partially tryptic + $pI^c$	468	5526
	forced tryptic	362	3875
	partially tryptic <sup>b</sup>	369	4223
	fully tryptic <sup>b</sup>	343	2531
	partially tryptic + $pI^c$	385	3909

<sup>a</sup> Results shown are obtained at  $P$  score or  $\Delta Cn$  values which give an estimated false positive rate of 1% using any given filtering method, as described in the text. <sup>b</sup> These results were obtained using no enzyme specification, and then filtering based on tryptic cleavage state as described in the text. <sup>c</sup> These results were obtained using no enzyme specification, then filtering for partially tryptic peptide sequences, followed by filtering by peptide  $pI$ , as described in the text.

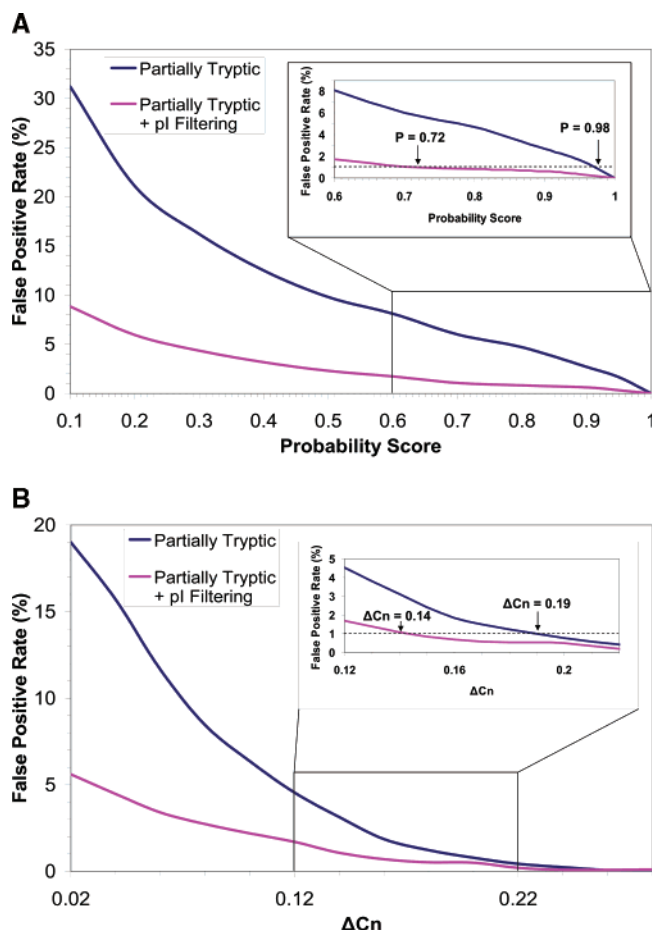
has been demonstrated that saliva contains biologically active peptide fragments derived from proteolysis other than trypsin.<sup>26,28</sup> Curve III in Figure 2A shows the results when filtering the data and only retaining matches to peptide sequences which are tryptic at both termini (i.e., fully tryptic peptide sequences). This filtering significantly decreases the false positive rate across all  $P$  score values, with the drawback that it eliminates correct

matches to partially tryptic peptide sequences which may be a significant portion of the total, as is the case for this sample (Table 1).

Figure 2B is a plot of the false positive rate against  $\Delta Cn$  when using traditional Sequest scoring parameters as described for Figure 1 above, and again compares the results using either the forced tryptic specification (curve I), or tryptic enzyme filtering (curves II and III). Clearly, the use of enzyme filtering drastically decreases the false positive rate when using  $\Delta Cn$  scoring. The number of proteins identified and peptide sequence matches using  $\Delta Cn$  scoring and these different filtering methods is shown in Table 1. Notably, when using enzyme filtering, the optimal  $\Delta Cn$  values necessary for keeping low false positive rates ( $\leq 1\%$ ) approach those which have been previously reported for 3D ion trap instruments (0.08<sup>9</sup> or 0.10<sup>4</sup>), although the optimal values for the LTQ data are still above these previously reported values. Furthermore, the  $P$  score strategy identifies more proteins than the use of  $\Delta Cn$  scoring, at least using this parameters specified here, providing increased support for the use of this type of probabilistic scoring routine in proteomic studies.

Importantly, Figure 2A and B also shows data comparing the results of data obtained from the LCQ to that obtained from the LTQ. For simplicity, we have plotted the false positive rates from the LCQ ion trap when searching with no enzyme specified and considering only fully tryptic peptide sequence matches (curve IV in Figure 2A and B). By comparison of curves III and IV in Figure 2A and B, it is apparent that despite the effectiveness in lowering false positive rates of specifying no enzyme in the database search parameters, followed by filtering for fully tryptic peptides, the rate of false positive matches is still higher for the LTQ than for the LCQ data when using either  $P$  score or  $\Delta Cn$  assignments to the sequence matches. This difference in false positive rates between the instruments was also the case when allowing partially tryptic peptide sequences as correct matches in the data filtering (data not shown). These results provide further evidence of the increased risk of false positives on the 2D ion trap instrument, and the need for more careful and stringent filtering of the MS/MS data derived from these instruments compared to traditional 3D ion trap instruments. It is also very clear from Figure 2A and B that using a forced tryptic constraint in the database search parameters dramatically increases the potential for false positive sequence matches when using the 2D ion trap data. Therefore, given the effectiveness in reducing the false positive rate, conducting sequence database searches with no enzyme specified, followed by filtering based on tryptic cleavage sequences, is a much preferable strategy for this type of MS/MS data.

**Addition of Peptide Physicochemical Information.** The use of physicochemical properties of peptides, introduced using separation strategies which have been termed information<sup>16</sup> or value-added<sup>29</sup> strategies, has also been shown to be helpful in shotgun proteomics to increase confidence in peptide sequence matches and reduce false positive peptide sequence matches. For example, it has been demonstrated<sup>30–32</sup> that the use of more accurate mass and charge state information of peptides measured by FTICR–MS can provide more confident sequence database search results. Other reports<sup>33–35</sup> have used retention time on the  $\mu$ LC reverse-phase column as a constraint in peptide sequence matches, demonstrating the ability to partially predict the elution time of peptide from reversed-phase columns and the use of this information in the peptide sequence identification process. Recently, the property of



**Figure 3.** Plot of the effect on the false positive peptide sequence match rate for the LTQ using no enzyme specification followed by partially tryptic enzyme filtering, and combined with peptide *pI* filtering, as described in the text. A. *P* score assignments; B. Sequest score assignments, varying the  $\Delta Cn$  score, with the following constant Xcorr threshold values:  $\geq 2.0$  (+1 ions);  $\geq 2.5$  (+2 ions); and  $\geq 3.8$  (+3 ions).

peptide isoelectric point (*pI*) introduced by preparative IEF of complex peptide mixtures by FFE<sup>16,18</sup> or immobilized pH gradient (IPG) gels<sup>10,36,37</sup> has been demonstrated to effectively reduce false positive matches, and also minimize false negative peptide sequence matches, which are sequence matches which are correct but do not score well enough to pass the peptide match scoring thresholds.<sup>10,16,18</sup> In a recent study of ours,<sup>16,18</sup> we showed that using peptide *pI* information of FFE fractionated peptides in combination with *P* score filtering we could relax the *P* score values while still maintaining a false positive rate below 1%, increasing the number of proteins identified at high confidence. For this particular study, because the complex peptide mixture of saliva peptides was fractionated based on peptide *pI* by FFE, we next investigated the use of peptide *pI* filtering, coupled with the trypsin enzyme cleavage strategy shown in Figure 2 as an approach to further maximize the number of high confidence protein identifications obtained using the LTQ instrument. Given the significant number of partially tryptic peptides contained in this saliva sample (Table 1), we applied the peptide *pI* filtering constraint to the peptide sequence matches originally filtered allowing for partially tryptic matches (curve II in Figure 2A and B), to maximize the number of peptide matches from this sample. In filtering these

data, we used a *pI* tolerance of 0.5 when comparing the *pI* of any given peptide sequence match and the measured *pI* of the FFE fraction being analyzed, which we have shown to be an optimal tolerance given the IEF resolution of the FFE instrument.<sup>18</sup> This filtering strategy is based on the assumption that the predicted peptide *pI* should closely correspond to the *pI* value of the FFE fraction containing the peptide.<sup>18</sup>

Figure 3 shows the results of using peptide *pI* information in conjunction with trypsin enzyme filtering to optimize high confidence protein identifications, with the false positive rate measured against both the *P* score value (Figure 3A) and the  $\Delta Cn$  value (Figure 3B). Clearly, the introduction of the *pI* constraint enables a relaxation of both the *P* score and  $\Delta Cn$  thresholds necessary to keep an estimated false positive rate of 1% or below. The net result of the introduction of the *pI* constraint is a significant increase in proteins identified and protein sequence coverage when using both *P* score values and  $\Delta Cn$  values (see Table 1). The additional proteins identified would otherwise be considered false negative matches which would not have satisfied the scoring criteria without the addition of the peptide *pI* information.<sup>10,18</sup> The combined use of trypsin enzyme filtering with peptide *pI* information as an added constraint therefore offers an optimal strategy for increasing high confident protein identifications and sequence coverage.

## Conclusions

The results of this study provide a note of caution to researchers using new 2D linear ion trap mass spectrometers for shotgun proteomics. The rapid acquisition of MS/MS data and consequently increased sensitivity of these instruments also comes with the tradeoff of increased potential for false positive peptide sequence matches. Therefore researchers should use care when analyzing these data and claiming protein identifications from these studies. This potential for false positive matches is especially apparent when forcing the sequence matches to match to tryptic cleavage sequences. We also have provided some effective solutions to reducing false positive matches and maximizing the number of high confidence protein identifications using 2D linear ion trap instruments. The use of trypsin enzyme filtering is a relatively easily implemented strategy to reduce false positives, which are significantly increased when using database search parameters forcing matches to only tryptic peptides. The only drawback with this filtering strategy is the increased time required (about 4 times longer for this study) for database searching of the large amount of MS/MS data obtained using the fast scanning 2D linear ion trap instrument, compared to database searching using a forced trypsin constraint, which may be a major limitation for researchers not having access to high throughput, parallel computing clusters. Furthermore, the addition of physiochemical constraints, such as peptide *pI*, combined with enzyme filtering is effective in further maximizing the number of high confidence protein identifications.

As a final note, given the dependence of scoring thresholds using traditional Sequest scoring on a variety of factors (e.g., search parameters, sequence database size, instrument used) which our present study and others<sup>11,12</sup> have demonstrated, proposing "optimal" scoring values for any given proteomic study is difficult, and we have purposely avoided proposing such parameters in this study given this fact. Our results support the use of more sophisticated scoring routines such as Peptide Prophet<sup>14</sup> or others,<sup>12</sup> which can maximize the



proteins identified at high confidence<sup>27</sup> and simplify the filtering of the data. However, the results shown using Peptide Prophet also indicates that the optimal thresholds for these advanced scoring routines are also dependent on database search parameters and the instrument used. Collectively, these findings point to the importance and utility of estimating false-positive rates for any given scoring criteria used to determine correct peptide sequence matches in large-scale proteomic studies. This false-positive estimate takes on new importance with data from the 2D linear ion trap, given the demonstrated increased potential for false positive matches when analyzing data from this instrumentation. It is our hope that commercial vendors would increasingly include with commercially available sequence database searching software an automated algorithm (e.g., reverse-database searching or other) which would allow for the routine estimation of false positive rates with protein identifications, enabling researchers to easily report this information with their proteomic datasets. Such implementation would help to address ongoing standardization efforts for the presentation and comparison of mass spectrometry-based proteomics data,<sup>38</sup> and aid researchers in the publication of only high quality data.

**Acknowledgment.** This work was supported in part by a grant from the Minnesota Medical Foundation. We thank Dr. Nelson L. Rhodus at the School of Dentistry of the University of Minnesota for human saliva sample, the Center for Mass Spectrometry and Proteomics at the University of Minnesota for instrumental resources, and the Minnesota Supercomputing Institute for computational support and maintenance of the Sequest cluster. T.J.G. also thanks the Eli Lilly and Company for financial support.

## References

- (1) Aebersold, R.; Mann, M. *Nature* **2003**, *422*, 198–207.
- (2) Yates, J. R., 3rd *Annu. Rev. Biophys. Biomol. Struct.* **2004**, *33*, 297–316.
- (3) Wolters, D. A.; Washburn, M. P.; Yates, J. R., 3rd *Anal. Chem.* **2001**, *73*, 5683–5690.
- (4) Washburn M. P.; Wolters, D. A.; Yates, J. R., 3rd *Nat. Biotechnol.* **2001**, *19*, 242–247.
- (5) Schwartz, J. C.; Senko, M. W.; Syka, J. E. *J. Am. Soc. Mass Spectrom.* **2002**, *13*, 659–669.
- (6) Mayya V.; Rezaul, K.; Cong, Y. S.; Han, D. *Mol. Cell. Proteomics* **2005**, *4*, 214–223.
- (7) Eng, J.; McCormack, A. L.; Yates, J. R., 3rd *J. Am. Soc. Mass Spectrom.* **1994**, *5*, 976–989.
- (8) Perkins, D. N.; Pappin, D. J.; Creasy, D. M.; Cottrell, J. S. *Electrophoresis* **1999**, *20*, 3551–3567.
- (9) Peng, J.; Elias, J. E.; Thoreen, C. C.; Licklider, L. J.; Gygi, S. P. *J. Proteome Res.* **2003**, *2*, 43–50.
- (10) Cargile, B. J.; Bundy, J. L.; Stephenson, J. L., Jr. *J. Proteome Res.* **2004**, *3*, 1082–1085.
- (11) Sadygov, R. G.; Yates, J. R., 3rd *Anal. Chem.* **2003**, *75*, 3792–3798.
- (12) Sadygov, R. G.; Liu, H.; Yates, J. R., 3rd *Anal. Chem.* **2004**, *76*, 1664–1671.
- (13) Bafna, V.; Edwards, N. *Bioinformatics* **2001**, *17 Suppl 1*, S13–21.
- (14) Keller, A.; Nesvizhskii, A. I.; Kolker, E.; Aebersold, R. *Anal. Chem.* **2002**, *74*, 5383–5392.
- (15) Krivankova L.; Bocek P. *Electrophoresis* **1998**, *19*, 1064–1074.
- (16) Xie H.; Rhodus, N. L.; Griffin, R. J.; Carlis, J. V.; Griffin, T. J. *Mol. Cell. Proteomics* **2005**, *4*, 1826–1830.
- (17) Rhodus, N. L.; Cheng, B.; Myers, S.; Bowels, W.; Ho, V.; Ondrey, F. *Clin. Immunol.* **2005**, *114*, 278–283.
- (18) Xie, H.; Bandhakavi, S.; Griffin, T. J. *Anal. Chem.* **2005**, *77*, 3198–3207.
- (19) Wilmarth, P. A.; Riviere, M. A.; Rustvold D. L.; Lauten, J. D.; Madden, T. E.; David, L. L. *J. Proteome Res.* **2004**, *3*, 1017–1134.
- (20) Han, D. K.; Eng, J.; Zhou, H.; Aebersold, R. *Nat. Biotechnol.* **2001**, *19*, 946–951.
- (21) Shimura, K.; Kamiya, K.; Matsumoto, H.; Kasai, K. *Anal. Chem.* **2002**, *74*, 1046–1053.
- (22) Yao, Y.; Berg, E. A.; Costello, C. E.; Troxler, R. F.; Oppenheim, F. G. *J. Biol. Chem.* **2003**, *278*, 5300–5308.
- (23) Vitorino, R.; Lobo, M. J.; Ferrer-Correia, A. J.; Dubin, J. R.; Tomer, K. B.; Domingues, P. M.; Amado, F. M. *Proteomics* **2004**, *4*, 1109–1115.
- (24) Ghafouri, B.; Tagesson, C.; Lindahl, M. *Proteomics* **2003**, *3*, 1003–1015.
- (25) Hardt, M.; Thomas, L. R.; Dixon, S. E.; Newport, G.; Agabian, N.; Prakobphol, A.; Hall, S. C.; Witkowska, H. E.; Fisher, S. J. *Biochemistry*, **2005**, *44*, 2885–2899.
- (26) Hu, S.; Xie, Y.; Ramachandran, P.; Ogorzalek Loo, R. R.; Li, Y.; Loo, J. A.; Wong, D. T. *Proteomics* **2005**, *5*, 1714–1728.
- (27) Kapp, E. A.; Schutz, F.; Connolly, L. M.; Chakel, J. A.; Meza, J. E.; Miller, C. A.; Fenyo, D.; Eng, J. K.; Adkins, J. N.; Omenn, G. S.; Simpson, R. J. *Proteomics* **2005**, *5*, 3475–3490.
- (28) Hardt, M.; Witkowska E.; Webb, S.; Thomas, L. R.; Dixon S. E.; Hall, S. C.; Fisher, S. J. *Anal. Chem.* **2005**, *77*, 4947–4954.
- (29) Heller, M.; Ye, M.; Michel, P. E.; Morier, P.; Stalder, D.; Jünger, M. A.; Aebersold, R.; Reymond, F.; Rossier, J. S. *J. Proteome Res.* **2005**, *4*, 2273–2282.
- (30) Qian, W. J.; Camp, D. G., II; Smith, R. D. *Expert Rev. Proteomics* **2004**, *1*, 87–95.
- (31) Ferguson, P. L.; Smith, R. D. *Annu. Rev. Biophys. Biomol. Struct.* **2003**, *32*, 399–424.
- (32) Lipton, M. S.; Pasa-Tolic, L.; Anderson, G. A.; Anderson, D. J.; Auberry, D. L.; Battista, J. R.; Daily, M. J.; Fredrickson, J.; Hixson, K. K.; Kostandarites, H.; Masselon, C.; Markillie, L. M.; Moore, R. J.; Romine, M. F.; Shen, Y.; Strittmatter, E.; Tolic, N.; Udseth, H. R.; Venkateswaran, A.; Wong, K. K.; Zhao, R.; Smith, R. D. *PNAS* **2002**, *99*, 11049–11054.
- (33) Petritis, K.; Kangas, L. J.; Ferguson, P. L.; Anderson, G. A.; Pasa-Tolic, L.; Lipton, M. S.; Auberry, K. J.; Strittmatter, E. F.; Shen, Y.; Zhao, R.; Smith, R. D. *Anal. Chem.* **2003**, *75*, 1039–1048.
- (34) Palmblad, M.; Ramstrom, M.; Markides, K. E.; Hakansson, P.; Bergquist, J. *Anal. Chem.* **2002**, *74*, 5826–5850.
- (35) Shen, Y.; Zhao, R.; Belov, M. E.; Conrads, T. P.; Anderson, G. A.; Tang, K.; Pasa-Tolic, L.; Veenstra, T. D.; Lipton, M. S.; Udseth, H. R.; Smith, R. D. *Anal. Chem.* **2001**, *73*, 1766–1775.
- (36) Cargile, B. J.; Tally, D. L.; Stephenson, J. L., Jr. *Electrophoresis* **2004**, *25*, 936–945.
- (37) Cargile, B. J.; Bundy, J. L.; Freeman, T. W.; Stephenson, J. L., Jr. *J. Proteome Res.* **2004**, *3*, 112–119.
- (38) Carr, S.; Aebersold, R.; Baldwin, M.; Burlingame, A.; Clauser, K.; Nesvizhskii, A. *Mol. Cell. Proteomics* **2004**, *3*, 531–533.

PR050472I