

Article

## MSblender: a probabilistic approach for integrating peptide identifications from multiple database search engines

Taejoon Kwon, Hyungwon Choi, Christine Vogel, Alexey I. Nesvizhskii, and Edward M. Marcotte

*J. Proteome Res.*, **Just Accepted Manuscript** • Publication Date (Web): 13 April 2011

Downloaded from <http://pubs.acs.org> on April 13, 2011

### Just Accepted

“Just Accepted” manuscripts have been peer-reviewed and accepted for publication. They are posted online prior to technical editing, formatting for publication and author proofing. The American Chemical Society provides “Just Accepted” as a free service to the research community to expedite the dissemination of scientific material as soon as possible after acceptance. “Just Accepted” manuscripts appear in full in PDF format accompanied by an HTML abstract. “Just Accepted” manuscripts have been fully peer reviewed, but should not be considered the official version of record. They are accessible to all readers and citable by the Digital Object Identifier (DOI®). “Just Accepted” is an optional service offered to authors. Therefore, the “Just Accepted” Web site may not include all articles that will be published in the journal. After a manuscript is technically edited and formatted, it will be removed from the “Just Accepted” Web site and published as an ASAP article. Note that technical editing may introduce minor changes to the manuscript text and/or graphics which could affect content, and all legal disclaimers and ethical guidelines that apply to the journal pertain. ACS cannot be held responsible for errors or consequences arising from the use of information contained in these “Just Accepted” manuscripts.

1  
2  
3 **MSblender: a probabilistic approach for integrating peptide**  
4 **identifications from multiple database search engines**  
5  
6  
7  
8  
9

10 Taejoon Kwon<sup>1\*</sup>, Hyungwon Choi<sup>2\*</sup>, Christine Vogel<sup>1,3</sup>, Alexey I. Nesvizhskii<sup>2,4‡</sup>, Edward M.  
11 Marcotte<sup>1,5‡</sup>  
12  
13  
14  
15  
16  
17

18 1. Center for Systems and Synthetic Biology, Institute for Cellular and Molecular Biology,  
19 University of Texas at Austin, Austin, Texas, USA  
20  
21

22 2. Department of Epidemiology and Public Health, Yong Loo Lin School of Medicine,  
23 National University of Singapore, Singapore  
24  
25

26 3. Center for Genomics and Systems Biology, Department of Biology, New York University,  
27 New York, USA  
28  
29

30 4. Department of Pathology and Center for Computational Medicine and Bioinformatics,  
31 University of Michigan, Ann Arbor, Michigan, USA  
32  
33

34 5. Department of Chemistry and Biochemistry, University of Texas at Austin, Austin, Texas,  
35 USA  
36  
37  
38  
39  
40

41 \* These authors contributed equally to this work.  
42  
43

44 ‡ To whom all correspondence should be addressed. Email: nesvi@med.umich.edu,  
45 marcotte@icmb.utexas.edu  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Abstract

Shotgun proteomics using mass spectrometry is a powerful method for protein identification but suffers limited sensitivity in complex samples. Integrating peptide identifications from multiple database search engines is a promising strategy to increase the number of peptide identifications and reduce the volume of unassigned tandem mass spectra. Existing methods pool statistical significance scores such as  $p$ -values or posterior probabilities of peptide-spectrum matches (PSMs) from multiple search engines after high scoring peptides have been assigned to spectra, but these methods lack reliable control of identification error rates as data are integrated from different search engines. We developed a statistically coherent method for integrative analysis, termed MSblender. MSblender converts raw search scores from search engines into a probability score for all possible PSMs and properly accounts for the correlation between search scores. The method reliably estimates false discovery rates and identifies more PSMs than any single search engine at the same false discovery rate. Increased identifications increment spectral counts for all detected proteins and allow quantification of proteins that would not have been quantified by individual search engines. We also demonstrate that enhanced quantification contributes to improve sensitivity in differential expression analyses.

**Keywords:** integrative analysis, database search, peptide identification

**Abbreviations:** FDR - false discovery rate, LC - liquid chromatography, MS - mass spectrometry, PSM - peptide-spectrum match, ROC - receiver operating characteristic

## Introduction

Analyses of mass spectrometry-based shotgun proteomics data rely heavily upon computational algorithms for automating peptide identification via database searching. Database search engines assign each tandem mass spectrum to the best-scoring peptide sequence in the database based on scoring functions using spectral features<sup>1-7</sup>. Several different search engines are available today, and peptides identified with high confidence often show good consensus across different engines<sup>8</sup>. Nevertheless, many high-quality MS/MS spectra remain unassigned to peptide sequences or have scores below chosen confidence thresholds<sup>9</sup>. Moreover, some spectra may be assigned to different peptides by different search engines, which vary in their scoring schemes<sup>4, 10</sup>. Provided that these issues are properly addressed, pooling peptide identifications from multiple search engines is expected to improve peptide identifications and to leave fewer mass spectra without assignment to peptide sequences.

To date, a few computational approaches have been proposed for integrating database search results. Alves *et al.* proposed a calibration of  $p$ -values from multiple search engines into a meta-analytic  $p$ -value for each peptide<sup>8</sup>. Searle *et al.* proposed a Bayes approach to adjust probability scores computed in individual search engines based on the agreement between search engines, in which the largest adjusted probability is taken as the final score for each peptide<sup>11</sup>. Although these methods allow for more efficient use of available data, the integration of search results still has room for further development. First, the number of peptide-spectrum matches (PSMs) identified in some but not all search engines grows at combinatorial rates as more search engines are considered for integration, and the scores must be properly calibrated for the PSMs identified by individual search engines to control the overall identification error rates in a unified manner. Second, since some search engines only report the best matching peptide sequence for each spectrum, potential matches to lower-ranking peptides are ignored in the report even if individual scores for those secondary matches are nearly as good as the best match score and thus

1  
2  
3 are likely true hits. If data are integrated from different search engines, one must include  
4 lower-ranking PSMs from every search engine and recalibrate the scores into a unified score  
5 as was done in Searle *et al*<sup>11</sup>. The strategy of integrating data after the selection of high  
6 confidence PSMs (*i.e.*, leaving out lower-ranking scores) may lead to inaccurate estimation  
7 of integrative probability scores unless search engines are sufficiently homogeneous<sup>12-14</sup>.  
8  
9

10  
11  
12 To address these issues, we developed a unified probabilistic approach for the  
13 integrative analysis of unique PSMs, termed MSblender (Figure 1). We use probability  
14 mixture models for distinguishing correct and incorrect identifications. The score distributions  
15 across search engines are jointly modeled using multivariate distributions up to the number  
16 of observed dimensions to accommodate the correlation in raw search scores. Using this  
17 model, MSblender computes a unified posterior probability of correct identification for all  
18 PSMs identified by search engines. The conversion into posterior probabilities automatically  
19 calibrates PSM scores reported by individual search engines in two ways: (1) the likelihood  
20 is marginalized to the search engines identifying individual PSMs, and (2) prior probability is  
21 adjusted for different combinations of search engines. More importantly, MSblender pools  
22 raw search scores for every possible PSM and directly models the distribution for all listed  
23 scores from the beginning, so it is not necessary to revisit lower-ranking PSMs to account for  
24 the PSMs not agreed upon by all search engines.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41

42 We evaluate the performance of MSblender with respect to peptide identification and  
43 protein quantification by spectral counting using three independent datasets. First, we use a  
44 yeast dataset (Yeast YPD hereafter) to assess the sensitivity and specificity profile for *bona*  
45 *fide* identifications, where high-confidence identifications reproducibly reported in multiple  
46 published datasets can be used as a benchmark set. Next, we include a (Sigma) UPS2  
47 dataset featuring a simple mixture of 48 human proteins, where concentrations are known  
48 for all proteins and thus the accuracy in both identification and quantification can be  
49 evaluated. Lastly, we use a dataset (iPRG09) from an Association of Biomolecular Resource  
50 Facilities (ABRF) proteome informatics research group (iPRG) 2009 study consisting of two  
51 biological samples, in which proteins present in only one sample are known and thus the  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 influence of improved identifications can be evaluated by differential expression analysis.  
4  
5 Through these examples, we show that integrative analysis by MSblender increases the  
6  
7 number of identifications substantially with accurate estimation of low false discovery rate  
8  
9 (FDR), and it improves quantitative analysis of protein concentrations.  
10  
11

## 12 13 14 **Materials and Methods**

### 15 16 17 **Yeast YPD dataset**

18  
19 Yeast YPD is a yeast dataset from Ramakrishnan *et al.*<sup>15</sup>. Briefly, cell lysates were  
20  
21 harvested from *S. cerevisiae* BY4741 grown in rich medium (YPD) in log phase, digested  
22  
23 with trypsin and prepared for LC/LC-MS/MS analysis. We performed eight replicate LC-  
24  
25 MS/MS using four salt steps on an SCX column (ammonium chloride solutions of varying  
26  
27 molarity, namely 0, 15, 60, 900 mM or 0, 20, 100, 900 mM in a 5% acetonitrile, 0.1% formic  
28  
29 acid background), followed by reverse-phase chromatography on a C18 column and MS/MS  
30  
31 analysis on an LTQ-Orbitrap Classic (Thermo). 32 files were analyzed using *S. cerevisiae*  
32  
33 sequences from EnsEMBL version 50 and randomly shuffled sequences as decoy. The raw  
34  
35 dataset is available at [http://www.marcottelab.org/users/MSdata/Data\\_02/](http://www.marcottelab.org/users/MSdata/Data_02/).  
36  
37  
38  
39

### 40 41 **UPS2 dataset**

42  
43 The dataset comprises 48 human proteins mixed in concentrations covering six orders of  
44  
45 magnitude, from 0.5 fmol to 50,000 fmol (Sigma Aldrich). The sample was prepared as  
46  
47 described before<sup>15</sup> including cysteine alkylation, trypsin digestion and cleanup of the  
48  
49 resulting peptides. The sample was re-suspended in 50  $\mu$ l of buffer (95% H<sub>2</sub>O, 5%  
50  
51 acetonitrile, 0.1% formic acid) and ten samples of different dilutions were used for LC-  
52  
53 MS/MS analysis on an LTQ-Orbitrap Classic (Thermo) mass spectrometer in a 5 to 90%  
54  
55 acetonitrile gradient over four hours. Dilutions ranged from none to 1:30, with 10 $\mu$ l injected  
56  
57 per run. We used a sequence file downloaded from Sigma Aldrich website as the target  
58  
59  
60

1  
2  
3 database and a decoy database derived from their randomly shuffled protein sequences.

4  
5 The raw data are deposited at [http://www.marcottelab.org/users/MSdata/Data\\_13/](http://www.marcottelab.org/users/MSdata/Data_13/).

### 6 7 8 9 **iPRG09 dataset**

10 We used the ABRF iPRG 2009 study data downloaded from Tranche Proteome Commons.

11  
12 The data consist of two 1D gel separations of identical *Escherichia coli* cellular lysates  
13  
14 (called the 'yellow' and 'red' samples). In each sample, one segment of the separation gel  
15  
16 was cut out and discarded. The two discarded segments ('green' and 'blue') did not overlap  
17  
18 in their position in the two samples, thus the proteins in these segments would be identified  
19  
20 as differentially expressed proteins relative to the other sample. For each of the two samples,  
21  
22 five LC-MS/MS data files were available. To compare our results with the original study, we  
23  
24 used the same *E. coli* sequences as available from the ABRF website

25  
26 <http://www.abrf.org/index.cfm/group.show/ProteomicsInformaticsResearchGroup.53.htm>

27  
28 including reversed sequences as decoy instead of randomly shuffled sequences.

### 29 30 31 32 33 34 35 **Data processing**

36 We used the same database with target and decoy sequences for all individual runs with  
37  
38 four different search engines: SEQUEST, X!Tandem with k-score, InsPecT and MyriMatch.

39  
40 We used default parameters in each search engine wherever possible, assuming that  
41  
42 parameters had already been optimized for each scoring matrix. We allowed for up to 2  
43  
44 missed tryptic cleavages, and set static cysteine alkylation. More detailed information, such  
45  
46 as software version and modified parameters, is reported in Table 1. In the search results,  
47  
48 individual spectra may be reported multiple times in different PSMs, mainly because of  
49  
50 different charge state assignments. For example, MyriMatch reported two PSMs for every  
51  
52 spectrum with different charge estimates (+2 and +3 in our default setup). To guarantee  
53  
54 consistency across search engines, we selected the best scoring PSM per spectrum based  
55  
56 on the scores listed in Table 1. However, this is not a requirement and additional lower  
57  
58  
59  
60

1  
2  
3 ranking PSMs might also be allowed if the total number of PSMs from each search engine is  
4  
5 not significantly different.  
6  
7

### 8 9 **Statistical model**

10 The statistical approach in MSblender is a probability mixture model for score distributions of  
11  
12 correct and incorrect identifications, which has been widely used for scoring PSMs<sup>16</sup>. A novel  
13  
14 feature of MSblender is that the mixture component distributions are modeled as multivariate  
15  
16 distributions that appropriately account for the correlation between database search scores  
17  
18 (Figure 2).  
19  
20  
21

22  
23 Suppose that the database search was repeatedly performed using  $K$  independent  
24  
25 search engines, and  $M$  spectra were matched to peptide sequences by at least one search  
26  
27 engine. Let  $S_i=(S_{i1}, S_{i2}, \dots, S_{iK})$  denote the raw search score for spectrum  $i$ , where  $i=1, 2, \dots, M$ .  
28  
29 We assume that the raw scores are median centered and scaled by setting unit standard  
30  
31 deviation in all search engines. Note that some  $S_{ij}$  can be missing if the  $j$ -th search engine  
32  
33 does not report the same PSM. The joint probability density of search scores can be written  
34  
35 as  
36  
37

$$38 \quad g(S_i) = (1-\pi) g_0(S_i) + \pi g_1(S_i) \text{ for all } i$$

39  
40 where  $\pi$  is the proportion of spectra with correct peptide assignment in the data.  $g_0$  and  $g_1$   
41  
42 are the score distributions for PSMs in correct and incorrect identifications, respectively. We  
43  
44 refer to them as negative and positive mixture component distributions from here on. To  
45  
46 estimate these distributions, a sufficient number of PSMs must have scores across all  
47  
48 database search algorithms included. However, not every spectrum is assigned to the same  
49  
50 peptide sequence across all database search engines (Figure 1). Especially for decoy  
51  
52 sequences, a spectrum is rarely assigned to the same decoy peptide by two or more search  
53  
54 engines due to the random nature of incorrect peptide matches. It follows that the negative  
55  
56 component  $g_0$  cannot be specified as a fully multivariate distribution due to the lack of usable  
57  
58 data for estimation, and thus we assume  $g_0(S_i) = \prod_{n=1, \dots, K} g_{0n}(S_{in})$ , *i.e.* scores from different  
59  
60



1  
2  
3 search engines are conditionally independent for incorrect identifications. Furthermore, it is  
4  
5 natural to assume different prior weights for true and false identifications when PSMs are  
6  
7 identified in more search engines than others, i.e. frequent identification implies high prior  
8  
9 belief of correct identification. Hence we vary the weight parameter  $\pi$  by each combination of  
10  
11 search engines, as many as  $2^K-1$ .  
12  
13

14 To provide flexibility for accommodating variable shapes of score distributions, we  
15  
16 allowed the mixture components  $g_0$  and  $g_1$  to be expressed as mixtures of multivariate  
17  
18 Gaussian distributions themselves ( $g_0$  with a diagonal covariance matrix), where the number  
19  
20 of subcomponents must be pre-specified by the user. Specifically,  
21  
22

$$23 \quad g_1(S_i) = \sum_{c=1, \dots, C} \lambda_c \text{MVN}_K(S_i; \mathbf{m}_c, \mathbf{V}_c)$$

24  
25 where  $c$  is the number of subcomponents for the positive component distribution  $g_1$ ,  $\text{MVN}_K$   
26  
27 stands for  $K$ -dimensional multivariate normal distribution, and  $\mathbf{m}_c$  and  $\mathbf{V}_c$  are the mean vector  
28  
29 and the covariance matrix for the subcomponent distribution  $c$  with a respective mixing  
30  
31 proportion  $\lambda_c$  (such that  $\sum_{c=1, \dots, C} \lambda_c = 1$ ). In a typical run, we specified two subcomponents by  
32  
33 default ( $c=2$ ).  
34  
35

36 In the case of  $g_0$ , the marginal negative component distribution  $g_{0n}$  of an individual  
37  
38 search engine  $n$  is expressed as a mixture of univariate Gaussian distributions to allow for  
39  
40 the same flexibility as in the positive component, i.e.  
41  
42

$$43 \quad g_0(S_i) = \prod_{n=1, \dots, K} \{ \sum_{c=1, \dots, C} \delta_{cn} N(S_{in}; m_{cn}, V_{cn}) \}$$

44  
45 where  $N$  denotes univariate normal distribution. Mean and variance parameters ( $m_{cn}$ ,  $V_{cn}$ ) as  
46  
47 well as the mixing proportion(s)  $\pi$  are estimated using the EM algorithm<sup>17</sup>, where the spectra  
48  
49 assigned to decoy peptides are treated as known incorrect PSMs, rendering the mixture  
50  
51 model semi-supervised<sup>18</sup>. Once we estimate the positive and negative distributions in the  
52  
53 score distribution, we compute the posterior probability of correct identification for each PSM  
54  
55 by Bayes' rule:  
56  
57

$$58 \quad p_i = P(\text{Correct} | S_i) = \pi g_1(S_i) / g(S_i)$$

1  
2  
3 where  $g(S_i) = (1-\pi) g_0(S_i) + \pi g_1(S_i)$  for every spectrum  $i$ , and  $\pi$  varies by search engine  
4 combinations. Recall that  $S_i$  is a fully  $K$ -dimensional vector without missing data only if PSM  $i$   
5 is observed in all search engines. For a spectrum with scores from fewer than  $K$  database  
6 search engines, we compute the probability using the marginal distributions of observed  
7 scores only. After computing probabilities, the FDR at a probability threshold  $p^*$  can be  
8 estimated by  $\sum_{i \in S^*} (1 - p_i) / |S^*|$ , where  $S^*$  is the set of PSMs such that  $p_i \geq p^*$  and  $|A|$  is the  
9 size of a set  $A$ .  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20

### 21 **Benchmark data**

22  
23 In the Yeast YPD dataset, we compared the MSblender results to the proteins observed in  
24 previously published large-scale data under the same condition (the entire cellular lysate  
25 during logarithmic growth in rich medium). This list of proteins was prepared from 4 MS-  
26 based proteomics datasets and 3 non-MS-based datasets (see Ramakrishnan *et al.*<sup>15</sup> and  
27 [http://www.marcottelab.org/MSdata/gold\\_yeast.html](http://www.marcottelab.org/MSdata/gold_yeast.html)). We used the list of 4,265 proteins  
28 observed in either two or more MS-datasets or any of non-MS-datasets as benchmark list.  
29  
30  
31  
32  
33  
34  
35  
36  
37

### 38 **Differential expression analysis by QSPEC**

39  
40 We applied a statistical method for selecting differentially expressed proteins based on  
41 spectral counts, termed QSPEC<sup>19</sup>, to the iPRG09 dataset analyzed by individual search  
42 engines and MSblender. QSPEC computes the odds (Bayes factors) of differential  
43 expression for individual proteins and reports log scaled odds multiplied by the sign  
44 determined by the direction of changes as the summary statistic. These quantities are used  
45 to estimate local fdr and FDR using nonparametric empirical Bayes methods<sup>20</sup>. We  
46 evaluated the sensitivity profile at various thresholds. We constructed receiver operating  
47 characteristic (ROC)-like curves using the benchmark set provided by the ABRF iPRG 2009  
48 study committee. The 'blue' and 'green' segments contain positive sets of enriched proteins  
49 in the 'red' and 'yellow' data respectively. In the ROC plot, the horizontal coordinate  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 corresponds to the number of proteins not included in the positive set, representing the false  
4  
5 positive hits. Likewise, the vertical coordinate corresponds to the number of proteins  
6  
7 included in the positive set, representing the sensitivity of detection.  
8  
9

## 10 11 **Software Availability**

12  
13  
14 The source for running MSblender can be downloaded from the URL  
15  
16 <http://www.marcottelab.org/index.php/MSblender>.  
17  
18

## 19 20 **Results and Discussion**

### 21 22 **Estimation of FDR**

23  
24  
25 The fundamental challenge in data integration is accurate estimation of error rates such as  
26  
27 FDR. This is particularly difficult when search engines are heterogeneous, *i.e.* conflicting  
28  
29 PSMs occur frequently between search engines, and search scores are not in good  
30  
31 agreement. Figure 3 plots the FDR estimated by MSblender (see Materials and Methods)  
32  
33 against decoy-based FDR estimates in all three datasets. We estimated decoy-based FDR  
34  
35 by labeling one half of the decoy PSMs as non-decoy PSMs and measuring their recurrence  
36  
37 in the MSblender results (with proper scaling). Overall, the two estimates show good  
38  
39 agreement in critical regions, *i.e.* where the error rate is low, in all datasets, with a trend of  
40  
41 underestimation of FDR against decoys in UPS2 and Yeast YPD datasets. There was no  
42  
43 evidence of such underestimation against decoys in iPRG09 datasets, particularly in the low  
44  
45 error rate area. A possible explanation for the underestimation is that more consistent decoy  
46  
47 PSMs were identified by multiple search engines in UPS2 and YPD datasets than in iPRG09  
48  
49 dataset. Since MSblender assigns higher prior weights for PSMs identified in more search  
50  
51 engines than for PSMs identified in only one engine, many multi-engine (borderline scoring)  
52  
53 decoy PSMs were assigned high probability.  
54  
55  
56

57  
58 To see this from a comparative angle, we first examined the union method, which  
59  
60 selects PSMs in individual search engines at fixed FDRs and merging them. In Table 2 (FDR

1  
2  
3 0.5%), rows for MSblender and union show that the latter approach consistently includes  
4  
5 more decoy PSMs than the former, leading to underestimation of error rates (actual error  
6  
7 rate by decoy count is higher than the target 0.5%). A more sophisticated union approach  
8  
9 with probability adjustments based on the search score agreement by Searle *et al*<sup>11</sup>  
10  
11 improves the accuracy of FDR estimates in the first two datasets but not in iPRG09 datasets  
12  
13 (Supplement Figure 2; see below). In addition to estimated FDR, Table 3 shows that  
14  
15 MSblender achieves a good trade-off between recovering more true targets (gold standard  
16  
17 reference identification in Yeast YPD dataset; see Materials and Methods) and identifying  
18  
19 false-positive proteins at FDR 1%, whereas union would have incurred higher error rates  
20  
21 than MSblender to achieve similar sensitivity. However, there was no significant  
22  
23 improvement at FDR 0.5%; union and MSblender reported nearly identical results, implying  
24  
25 that at extremely low error cutoffs, a sophisticated statistical model, such as is used in  
26  
27 MSblender, is not necessarily helpful.  
28  
29  
30

31  
32 The findings above show that choosing thresholds for each search engine before  
33  
34 forming the union is non-trivial as merging data filtered at a fixed FDR in individual search  
35  
36 engines results in the post-integration FDR being higher than the target FDR. In addition,  
37  
38 there may not exist a unique solution to control the composite identification error rate since  
39  
40 the combined error is not a simple function of the individual error rates. In this situation, FDR  
41  
42 control by MSblender based on multivariate mixture model can be helpful despite its  
43  
44 underestimation of error in extremely high confidence regions.  
45  
46  
47  
48

### 49 **Sensitivity of identification**

50  
51 In addition to the accuracy of FDR estimates in high confidence selections, we evaluated the  
52  
53 performance of MSblender in comparison to the individual search engines. Specifically, we  
54  
55 examined the number of identifications for both PSMs and proteins at fixed FDRs (0.5% and  
56  
57 1%), and summarized the result in Figure 4. In all three datasets, MSblender clearly  
58  
59 increases the number of identifications over individual search engines. For example,  
60  
X!Tandem yielded the highest number of PSMs amongst other search engines at FDR 0.5%

1  
2  
3 in the Yeast YPD dataset (Figure 4a). Table 2 shows that the dataset contains 240,781  
4 PSMs, and X!Tandem identified 74,244 PSMs among these (30%). In comparison,  
5 MSblender identified 99,814 PSMs (41%) at the same FDR, increasing the number of  
6 identified PSMs by 11% of total spectra. From these additional PSMs, 452 new proteins  
7 were identified by MSblender in comparison to X!Tandem in this dataset. In the UPS2  
8 dataset, the number of PSMs increased by ~20% in MSblender compared to the best  
9 performing search engine SEQUEST (Figure 4b), but the number of proteins increased only  
10 marginally (by 4 proteins) due to the low sample complexity (see Table 2). However, the  
11 additional PSMs helped to improve quantification by spectral counting (see below). In  
12 iPRG09 datasets (Figures 4c-4d), MSblender identified a substantial number of additional  
13 PSMs, e.g. roughly doubling the number of identified PSMs over SEQUEST, and many new  
14 proteins.  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

29 In general, MSblender consistently increased the number of identifications compared  
30 to individual search engines both when integrating all search engines and when integrating  
31 only a subset of the engines (Supplementary Figure 1). This result holds for both small and  
32 highly complex protein samples with thousands of proteins. Interestingly, MSblender was as  
33 good with some combinations of three search engines as with all four search engines,  
34 indicating that there exists saturation effect in terms of additional improvements by including  
35 additional search engines.  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45

### 46 **Comparison to agreement score adjustment**

47  
48 To compare MSblender to existing integrative methods, we implemented in-house code to  
49 carry out the procedure described in Searle *et al.*<sup>11</sup>, which we term agreement score  
50 adjustment (personal communication with the author). For each PSM, the agreement score  
51 was calculated following Searle *et al.*'s description, the probability scores from individual  
52 search engines were adjusted reflecting this agreement score, and the largest probability  
53 from all available search engines was taken as the final score for each PSM. If a search  
54 engine does not report the PSM, we considered it as a zero probability event.  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Supplementary Figure 2 shows that agreement score adjustment not only produces as many PSMs as MSblender at fixed FDRs in UPS2 dataset, but also provides more accurate FDR estimates. The performance for the two methods was also similar in Yeast YPD dataset. However, MSblender offers much more accurate FDR estimates and more PSMs in both iPRG09 datasets. This is partly attributable to the fact that MSblender performs very well with highly heterogeneous search engine results, i.e. search engines that largely disagree on PSM identifications, as is the case in iPRG09 but not in UPS2 and Yeast YPD. Supplementary Table 1 shows that both decoy and non-decoy PSMs common in two or more search engines appeared more frequently in Yeast YPD and UPS2 datasets than in iPRG09 datasets. In the iPRG09 data, the method of agreement score adjustment operates close to the union method as it takes the largest probability, but probability adjustment occurs less frequently in iPRG09 than in other datasets (Supplementary Figure 2 panels 3b-4b).

### Impact on label-free quantification by spectral counting

Enhanced PSM and protein identification by integrative analysis is expected to have a positive impact on spectral counting based protein quantification, in particular if additional PSMs contribute to total spectral counts per protein and thus refine the resolution of quantification. For instance, the difference between 1 spectrum and 2 spectra is less discernible than the difference between 10 spectra and 20 spectra. The increase in spectral counts is particularly important for low-abundance proteins, as they are often identified by one or a few PSMs only using a single search engine.

We first examined the number of proteins identified by individual search engines and MSblender across the range of spectral counts. In spectral counting, if a spectrum is matched to  $n$  different peptides at a given significance level, then we considered the PSM as  $1/n$  count for each and every peptide (split spectral counts). We also investigated not using this correction and adding one count to each matched peptide, but this procedure did not perform differently in the low FDR range (*not shown*). Supplementary Figures 3a-3d illustrate

1  
2  
3 the results, for the case of split spectral counts. When we plotted the number of identified  
4 proteins against the spectral counts normalized by their length (in  $\log_{10}$  scale), we found that  
5 not only the spectral count per protein increased in MSblender compared to individual  
6 searches, but also additional low-abundance proteins were identified. This observation  
7 implies that MSblender not only increased spectral counts for proteins identified by individual  
8 search engines, but also identified additional proteins in low spectral counts that individual  
9 search engines missed.

10  
11 To assess the accuracy of quantification, we examined the correlation between  
12 spectral counts of 48 proteins of known concentrations (UPS2 dataset). Figure 5 shows the  
13 results for the individual search engines and MSblender combining all of them at FDR 0.5%.  
14 As in Yeast YPD data, spectral counts were normalized by protein length, and spectral  
15 counts and known concentrations were rescaled to  $\log_{10}$ . Figures 5a-5e show that  
16 MSblender allows for quantification of the largest number of proteins (42 proteins out of 48)  
17 while providing a similar correlation between observed and expected protein concentrations  
18 as individual searches, measured by rank correlation coefficients (Figure 5f). Even though  
19 tens of thousands of additional PSMs have been used by MSblender for quantification,  
20 length-normalized spectral counts were linearly correlated with the known concentrations.  
21 Based on the evenly distributed increase in identifications, we expect that this trend should  
22 hold in samples of high complexity. When we relaxed the error criterion (FDR 1%), the  
23 number of quantified proteins also increased in individual search engines, but the correlation  
24 decreased. FDR cutoffs less stringent than 1% admitted PSMs with probability score as low  
25 as 0.5 or 0.6, allowing noisy PSMs to compromise quantification accuracy.

26  
27 To further demonstrate that increased PSM identification improves quantification, we  
28 used iPRG09 datasets to examine the sensitivity/specificity profile during differential  
29 expression analysis. To evaluate the performance consistently, we applied the differential  
30 expression analysis tool QSPEC (see Materials and Methods) to spectral count data  
31 reported by individual search engines and MSblender. Figure 6 shows the comparative  
32 performance in terms of the receiver-operating characteristic (ROC) for split and raw

1  
2  
3 spectral count data at FDR 0.5%. MSblender detected differentially expressed proteins with  
4  
5 the highest sensitivity at all fixed error rates for both types of spectral count data.  
6  
7

## 8 9 10 **Conclusion**

11  
12 This work presents an efficient probabilistic approach that substantially improves  
13 identification of PSMs at low error rates, reducing the volume of unassigned spectra in mass  
14 spectrometry-based shotgun proteomics experiments. Because scores are computed for all  
15 PSMs from the start, PSMs subject to search engine discrepancy are automatically  
16 considered for all possible peptide sequences and thus there is no need to trace back lower-  
17 ranking PSMs in individual searches. The final probability of correct identification does not  
18 have to rely on the probability from individual search engines which was calculated ignoring  
19 statistical dependence of raw search scores. The method can be applied to any number and  
20 combination of database search engines. Since the score is directly calculated from raw  
21 scores, the underlying statistical model allows a unified control of identification error rates.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

34 Integration of unique PSMs in MSblender provides justifiable grounds for more  
35 coherent quantification than the post-assignment integration methods, particularly if spectral  
36 counting is employed. Interestingly, we observed that the identification and quantification  
37 improved not only for low abundance proteins with MSblender, but also for high abundance  
38 proteins. This observation implies that MSblender substantially expanded the dynamic range  
39 of detectable protein concentrations without compromising quantification accuracy.  
40  
41  
42  
43  
44  
45  
46

47 For practical applications, we remind the readers that the fundamental challenge of  
48 integrative analysis is accurate estimation of identification errors, and it is important to  
49 understand the benefits and risks for error estimation when integrating different search  
50 engines of varying heterogeneity. MSblender delivers a tool to estimate correct integrated  
51 error rates even across heterogeneous datasets. In homogeneous datasets, *i.e.* where many  
52 search engines share PSMs including decoy PSMs, the performance of MSblender can be  
53 improved by modeling the negative component distribution in multivariate form as the  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 positive component distribution. In practice this is not feasible because the proportion of  
4 decoy PSMs occurring in two or more search engines is too small. We leave further  
5 improvements in statistical modeling to future work.  
6  
7  
8  
9

## 10 11 **Acknowledgments**

12 We thank Dr. Brian Searle for generously providing his code excerpt from Scaffold. This  
13 work was supported grants from the NIH, NSF, the Welch (F1515) & Packard Foundations  
14 (to E.M.M.), and NIH grants R01-GM094231 and R01-CA126239 (to A.I.N.).  
15  
16  
17  
18  
19  
20  
21  
22

## 23 **Supporting Information**

24 Supporting Information is available free of charge via the Internet at <http://pubs.acs.org>.  
25  
26  
27  
28  
29

30 **Supplementary Figure 1.** Comparison of identification results across different database  
31 search engines and MSblender integrating 2, 3, and 4 search engines. The counts of PSMs  
32 assigned to non-decoy sequences are plotted versus the estimated FDR. The figure shows  
33 that MSblender integration using combinations of three search engines performs  
34 comparable to MSblender with all four search engines.  
35  
36  
37  
38  
39  
40  
41  
42

43 **Supplementary Figure 2.** Comparison of MSblender and the agreement score method of  
44 Searle *et al.* with respect to the number of identifications (1a-4a) and the consistency  
45 between the estimated FDR and the decoy-based FDR (1b-4b). MSblender identifies more  
46 PSMs than the agreement score method except for the UPS2 dataset. The agreement  
47 method controls errors more accurately (closer to the diagonal) in the YPD and UPS2  
48 datasets; MSblender estimates errors more accurately for the iPRG09 dataset.  
49  
50  
51  
52  
53  
54  
55  
56  
57

58 **Supplementary Figure 3.** The number of identified proteins as a function of MS/MS spectral  
59 counts. Spectral counts per protein were normalized by protein length and rescaled to  $\log_{10}$ .  
60

1  
2  
3 The improvement to protein identification by MSblender is similar for proteins from all  
4 abundance ranges, and does not show strong bias e.g. towards high abundance proteins.  
5  
6  
7  
8

9  
10 **Supplementary Table 1.** Heterogeneity of PSM agreement across search engines before  
11 filtering.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Reference

1. Craig, R.; Beavis, R. C., TANDEM: matching proteins with tandem mass spectra. *Bioinformatics* **2004**, *20*, (9), 1466-7.
2. Eng, J. K.; McCormack, A. L.; Yates Iii, J. R., An approach to correlate tandem mass spectral data of peptides with amino acid sequences in a protein database. *Journal of the American Society for Mass Spectrometry* **1994**, *5*, (11), 976-989.
3. Geer, L. Y.; Markey, S. P.; Kowalak, J. A.; Wagner, L.; Xu, M.; Maynard, D. M.; Yang, X.; Shi, W.; Bryant, S. H., Open mass spectrometry search algorithm. *J Proteome Res* **2004**, *3*, (5), 958-64.
4. 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., An evaluation, comparison, and accurate benchmarking of several publicly available MS/MS search algorithms: sensitivity and specificity analysis. *Proteomics* **2005**, *5*, (13), 3475-90.
5. Perkins, D. N.; Pappin, D. J.; Creasy, D. M.; Cottrell, J. S., Probability-based protein identification by searching sequence databases using mass spectrometry data. *Electrophoresis* **1999**, *20*, (18), 3551-67.
6. Tabb, D. L.; Fernando, C. G.; Chambers, M. C., MyriMatch: highly accurate tandem mass spectral peptide identification by multivariate hypergeometric analysis. *J Proteome Res* **2007**, *6*, (2), 654-61.
7. Tanner, S.; Shu, H.; Frank, A.; Wang, L. C.; Zandi, E.; Mumby, M.; Pevzner, P. A.; Bafna, V., InsPecT: identification of posttranslationally modified peptides from tandem mass spectra. *Anal Chem* **2005**, *77*, (14), 4626-39.
8. Alves, G.; Wu, W. W.; Wang, G.; Shen, R. F.; Yu, Y. K., Enhancing peptide identification confidence by combining search methods. *J Proteome Res* **2008**, *7*, (8), 3102-13.
9. Ning, K.; Fermin, D.; Nesvizhskii, A. I., Computational analysis of unassigned high-quality MS/MS spectra in proteomic data sets. *Proteomics* **2010**, *10*, (14), 2712-8.
10. Sultana, T.; Jordan, R.; Lyons-Weiler, J., Optimization of the Use of Consensus Methods for the Detection and Putative Identification of Peptides via Mass Spectrometry Using Protein Standard Mixtures. *J Proteomics Bioinform* **2009**, *2*, (6), 262-273.
11. Searle, B. C.; Turner, M.; Nesvizhskii, A. I., Improving sensitivity by probabilistically combining results from multiple MS/MS search methodologies. *J Proteome Res* **2008**, *7*, (1), 245-53.
12. Lu, P.; Vogel, C.; Wang, R.; Yao, X.; Marcotte, E. M., Absolute protein expression profiling estimates the relative contributions of transcriptional and translational regulation. *Nat Biotechnol* **2007**, *25*, (1), 117-24.
13. Old, W. M.; Meyer-Arendt, K.; Aveline-Wolf, L.; Pierce, K. G.; Mendoza, A.; Sevinsky, J. R.; Resing, K. A.; Ahn, N. G., Comparison of label-free methods for quantifying human proteins by shotgun proteomics. *Mol Cell Proteomics* **2005**, *4*, (10), 1487-502.
14. Zybailov, B.; Mosley, A. L.; Sardi, M. E.; Coleman, M. K.; Florens, L.; Washburn, M. P., Statistical analysis of membrane proteome expression changes in *Saccharomyces cerevisiae*. *J Proteome Res* **2006**, *5*, (9), 2339-47.
15. Ramakrishnan, S. R.; Vogel, C.; Prince, J. T.; Li, Z.; Penalva, L. O.; Myers, M.; Marcotte, E. M.; Miranker, D. P.; Wang, R., Integrating shotgun proteomics and

1  
2  
3 mRNA expression data to improve protein identification. *Bioinformatics* **2009**, 25,  
4 (11), 1397-403.

5  
6 16. Keller, A.; Nesvizhskii, A. I.; Kolker, E.; Aebersold, R., Empirical statistical  
7 model to estimate the accuracy of peptide identifications made by MS/MS and  
8 database search. *Anal Chem* **2002**, 74, (20), 5383-92.

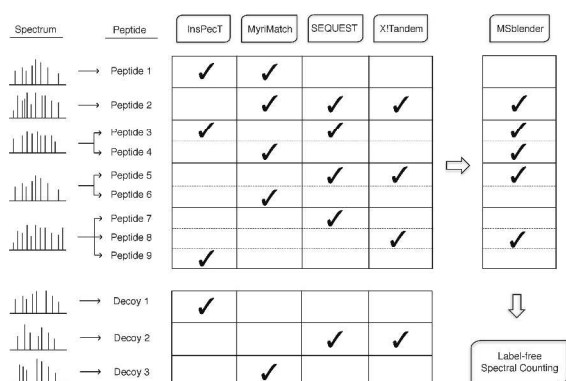
9  
10 17. Dempster, A. P.; Laird, N. M.; Rubin, D. B., Maximum likelihood from  
11 incomplete data via the EM algorithm. *J. Royal Statist. Soc.* **1977**, 39, (1), 1-38.

12  
13 18. Choi, H.; Nesvizhskii, A. I., Semisupervised model-based validation of peptide  
14 identifications in mass spectrometry-based proteomics. *J Proteome Res* **2008**, 7, (1),  
15 254-65.

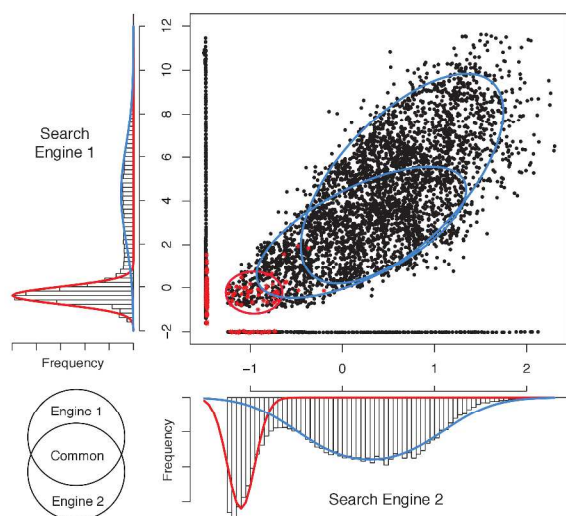
16  
17 19. Choi, H.; Fermin, D.; Nesvizhskii, A. I., Significance analysis of spectral count  
18 data in label-free shotgun proteomics. *Mol Cell Proteomics* **2008**, 7, (12), 2373-85.

19  
20 20. Efron, B.; Tibshirani, R.; Storey, J. D.; Tusher, T., Empirical Bayes analysis of  
21 a microarray experiment. *J Am Statist Assoc* **2001**, 96, 1151-1160.

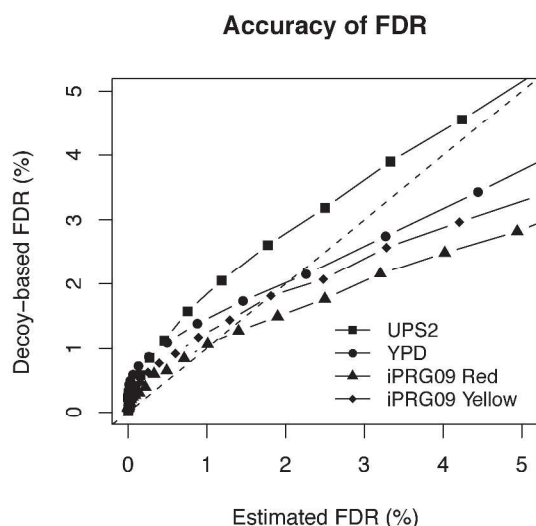
## Figures



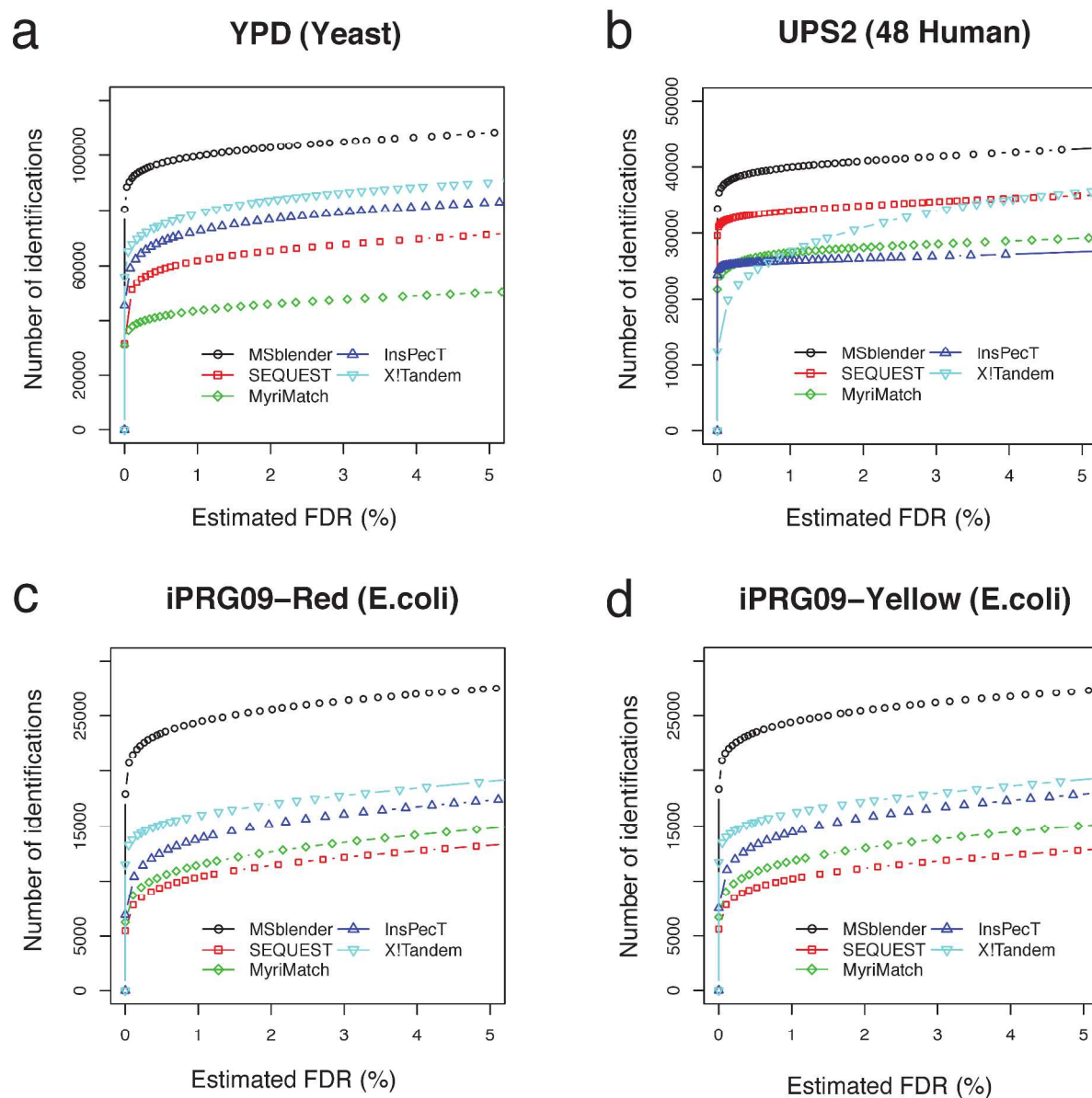
**Figure 1.** Schematic view of MSblender. Each spectrum is listed with all unique peptide assignments across the database search engines considered for integration. MSblender considers all different cases. PSMs may be found by some search engines, but not by others (peptides 1, 2). Some spectra may be matched to different peptides by different search engines (peptides 3 - 9), where each PSM is treated differently.



**Figure 2.** An example of the statistical model for integrating scores from two database search engines. The scatter plot shows three groups of data: PSMs with scores reported from both engines (dots following a diagonal line), and PSMs identified uniquely by either of the two search engines (dots in lines parallel to each axis). Red stars indicate the PSMs assigned to decoy sequences. The elliptical contours in the scatter plot and the curves in the histograms are the estimated distributions (blue: correct identification, red: incorrect identification).

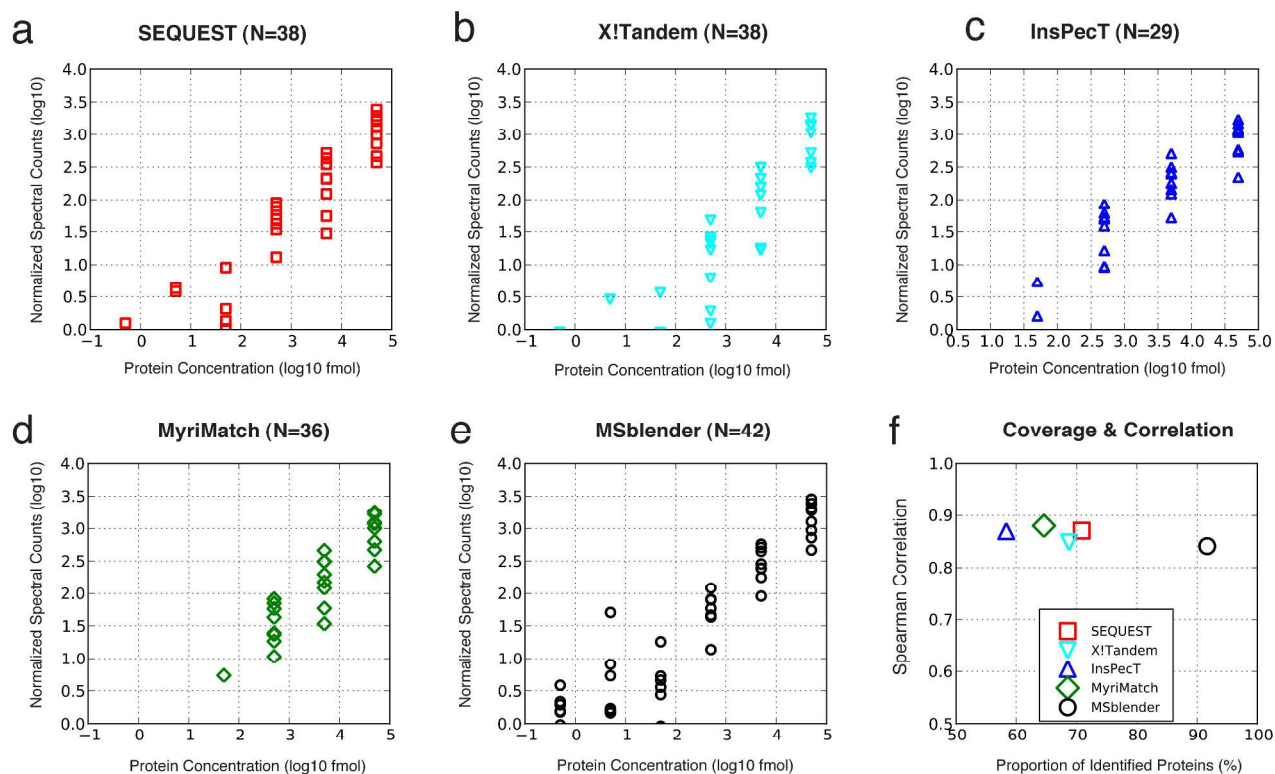


**Figure 3.** FDR estimated from posterior probabilities (Estimated FDR) against FDR estimated from decoy identifications (Decoy FDR). Estimated FDR is calculated by averaging PSM errors with a posterior probability threshold as described in Materials and Methods. Decoy-based FDR is calculated by recovery rate of decoy PSMs after labeling a half of decoy PSMs as target PSMs before running MSblender. Provided that decoys are truly random hits, the diagonal line indicates accurate FDR estimates.

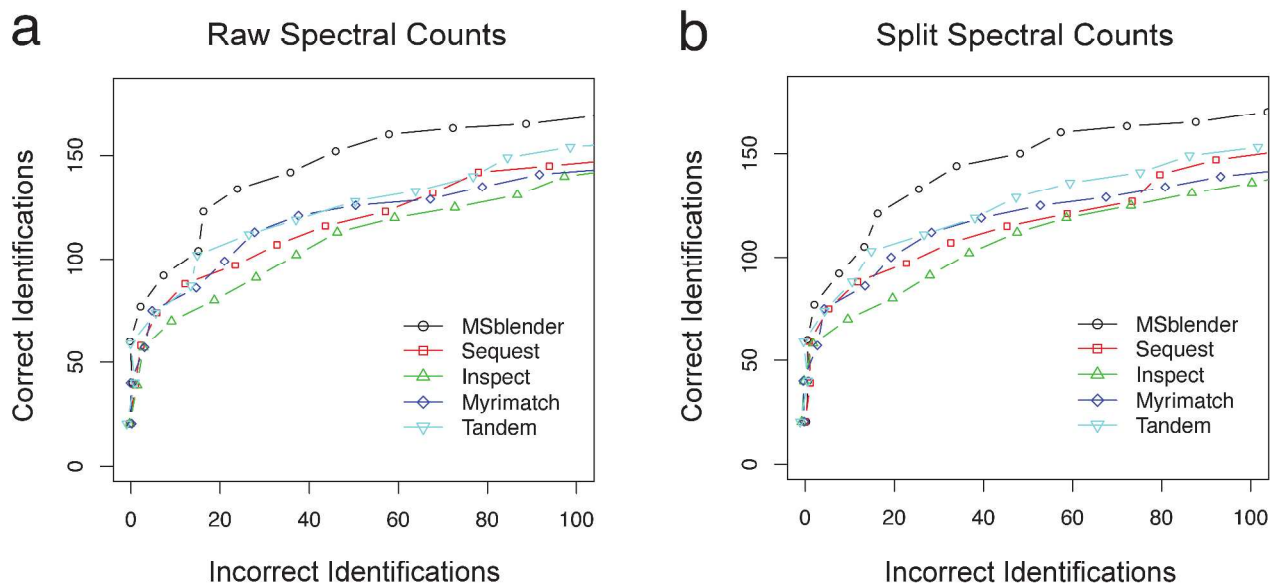


**Figure 4.** Comparison of identification results across different database search engines and MSblender integrating all search engines. The plots show the number of PSMs assigned to non-decoy sequences against estimated FDR. Although the performance of individual search engine varies depending on the dataset, it is clear that integrated MSblender can identify substantially more PSMs at the same FDR. More figures with different combination of search engine results are available in the Supplementary Figure 1.





**Figure 5.** Length normalized spectral counts against known protein concentrations in the UPS2 dataset (FDR 0.5%). Spectral counts and concentrations were rescaled to log<sub>10</sub>. **a-e:** Individual search engines and MSBlender. The number in subtitle parentheses reports the number of proteins identified. Across all protein concentrations, normalized spectral counts increase in MSBlender, showing that MSBlender improves the dynamic range of spectral counts. **f:** Spearman rank correlation coefficients ( $R_s$ ) against percentage of proteins identified, defined as the number of identified proteins divided by the total number of proteins known to exist (48). MSBlender improves the protein identification substantially (x-axis), while maintaining a correlation between observed and known protein concentrations similar to single search engine results.



**Figure 6.** Comparison of QSPEC analyses using the search results from individual search engines and MSblender integrating all search engines (FDR 0.5%). Correct and incorrect identifications were determined from proteins removed from each sample provided by the ABRF iPRG 2009 committee. See Materials and Methods for details in each samples and segments.

## Tables

**Table 1.** Summary of search engine parameters. Parameters not reported in this table were not changed from default values used by the search engine.

Name	Source	Version	Scores used in MSblender	Parameters
SEQUEST	Thermo Electron	Bioworks 3.3.1 SP1	Xcorr	Mass type: monoisotopic precursor and fragments Peptide tolerance: 25.0 ppm Fragment ion tolerance: 1.0 amu
X!Tandem (k-score)	REFERENCE	2009.10.01.1	E-value	Fragment monoisotopic mass error: 0.7 Parent monoisotopic mass error: 100 ppm Minimum peaks: 15 Minimum fragment m/z: 150
InsPecT	REFERENCE	20100331	MQscore	TagCount: 50 PMTolerance: 2.5
MyriMatch	REFERENCE	1.6.62 (2009-12-4)	Mvh	NumChargeStates: 3 UseAvgMassOfSequences: false

**Table 2.** Summary table of identification results by individual search engines and MSblender combining all search engine results. The entries were obtained at FDR 0.5%. In the PSM table, the rows referred to as 'k engines' indicate the number of PSMs identified with k search engines. In the same table, the numbers in the parentheses are the number of decoy PSMs identified at the same FDR (0.5%).

PSM	UPS2	Yeast YPD	iPRG09(Red)	iPRG09(Yellow)
Total MS/MS spectra observed	74,602	240,781	69,416	70,970
SEQUEST	32,651 (87)	57,955 (268)	9,524 (98)	9,492 (83)
X!Tandem	27,264 (210)	74,244 (332)	15,147 (117)	15,366 (112)
MyriMatch	26,262 (79)	41,179 (106)	9,706 (88)	9,134 (46)
InsPecT	25,618 (64)	69,341 (414)	12,691 (202)	13,295 (216)
Union	40,829 (434)	95,315 (1053)	21,764 (505)	21,684 (455)
MSblender	39,273 (336)	99,814 (1011)	23,580 (153)	23,717 (177)
1 engine	4,043 (190)	10,441 (100)	2,138 (38)	2,073 (52)
2 engines	7,389 (89)	16,861 (546)	3,768 (76)	3,878 (74)
3 engines	5,560 (35)	32,111 (203)	6,820 (24)	6,816 (21)
4 engines	22,202 (24)	38,257 (18)	10,830 (3)	10,826 (4)
Protein	UPS2	Yeast YPD	iPRG09(Red)	iPRG09(Yellow)
Total proteins	48	6,698	4,417	4,417
SEQUEST	38	1,391	757	749
X!Tandem	38	1,459	870	847
MyriMatch	36	1,241	722	657
InsPecT	29	1,527	877	902
Union	44	1,873	999	1,024
MSblender	42	1,911	1,185	1,147

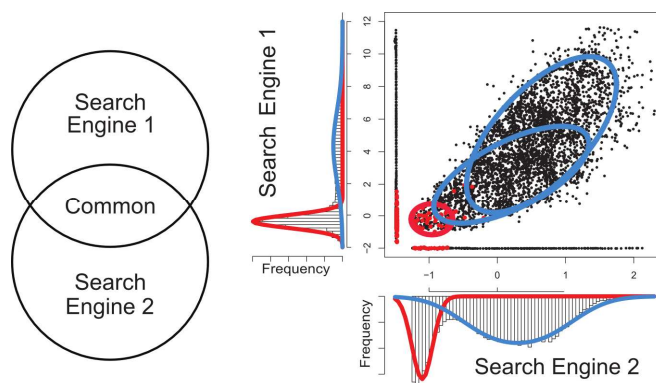
**Table 3.** Recovery of gold standard proteins in Yeast YPD dataset.

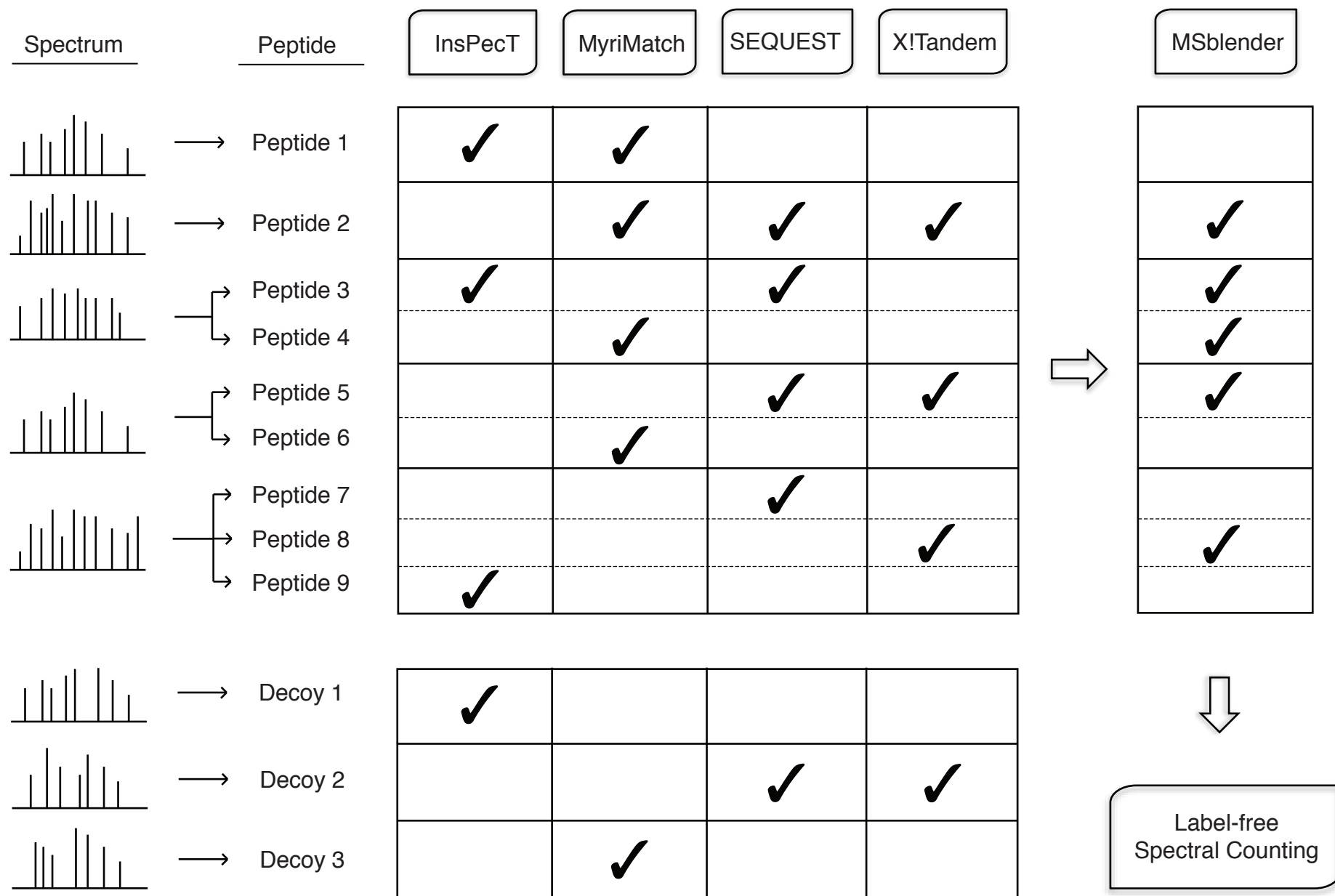
FDR 0.5%	All proteins	False proteins	True proteins	False/Total
Total	6698	2433	4265	
SEQUEST	1387	69	1318	4.97%
X!Tandem	1453	62	1268	4.26%
MyriMatch	1238	43	1195	3.47%
InsPecT	1519	99	1420	6.51%
Union	1899	161	1738	8.47%
MSblender	1864	153	1711	8.20%

FDR 1%	All Proteins	False proteins	True proteins	False/Total
Total	6698	2433	4265	
SEQUEST	1500	96	1404	6.40%
X!Tandem	1628	98	1530	6.01%
MyriMatch	1307	53	1254	4.05%
InsPecT	1662	134	1528	8.06%
Union	2218	252	1966	11.36%
MSblender	2038	203	1835	9.96%

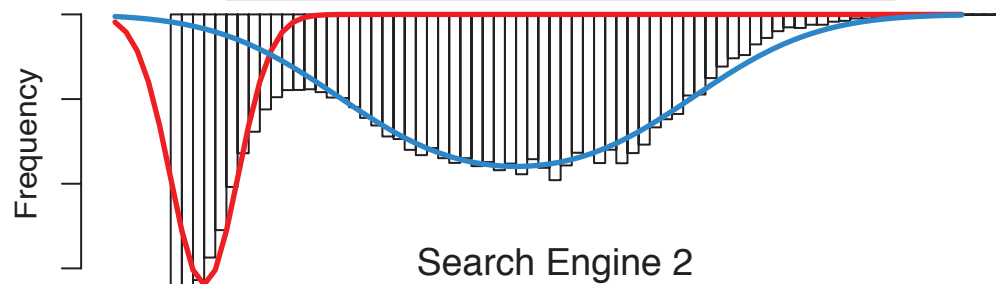
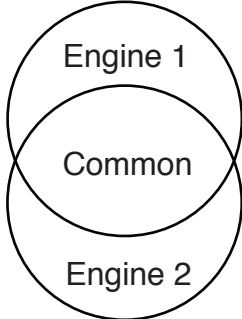
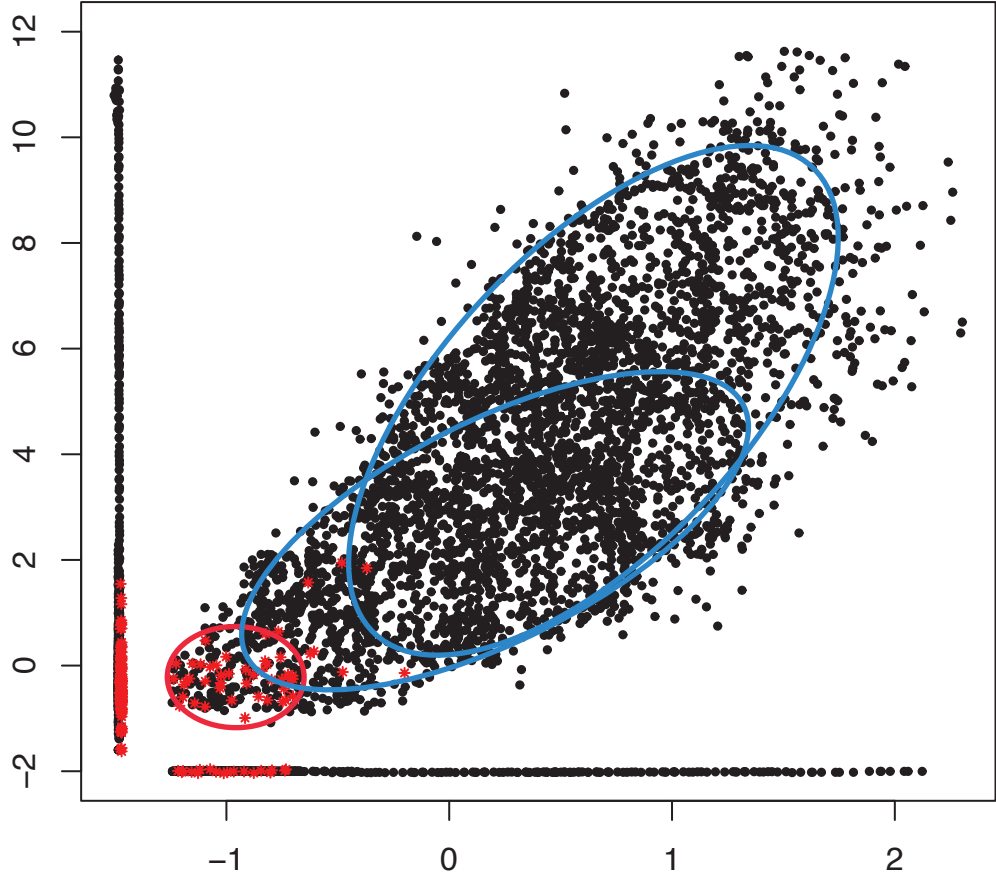
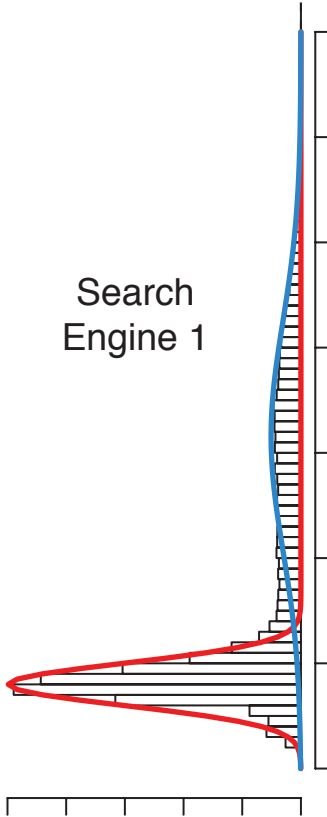
### Table of Contents Synopsis

We present a statistical method, termed MSblender, which integrates mass spectrometry based peptide identification results from multiple database search engines. MSblender models the joint score distributions of peptide-spectrum matches in multiple search engines, and performs classification of each match into correct or incorrect identifications. It increases the number of peptide identifications significantly at low false discovery rates. The improvement also leads to better quantification of low abundance proteins.



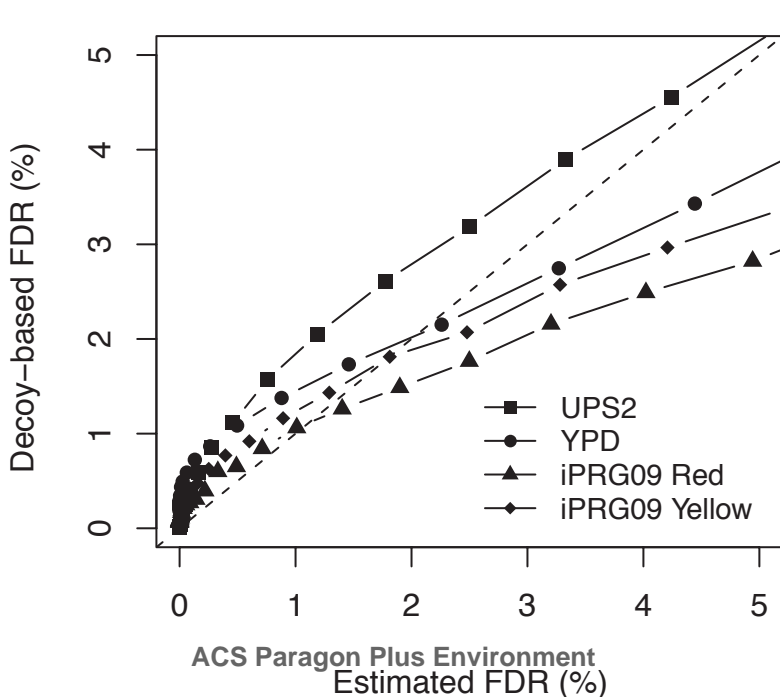


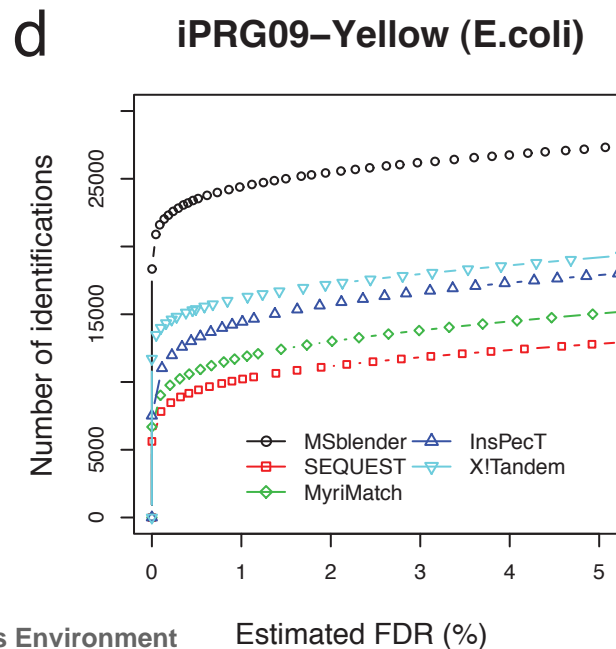
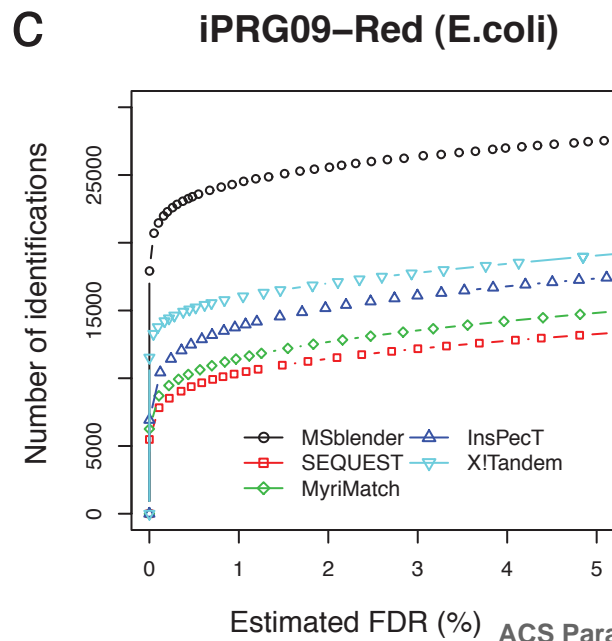
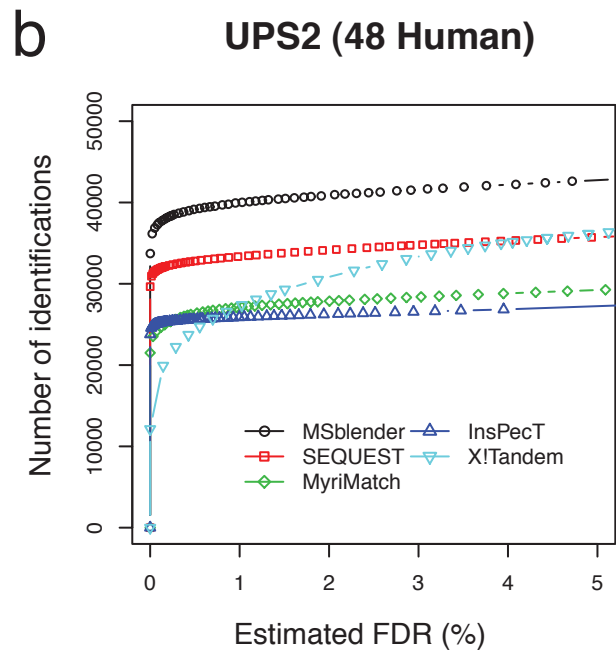
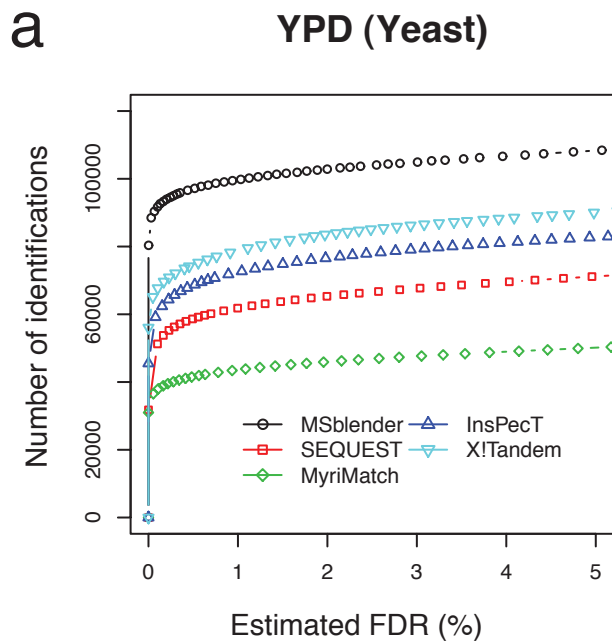
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

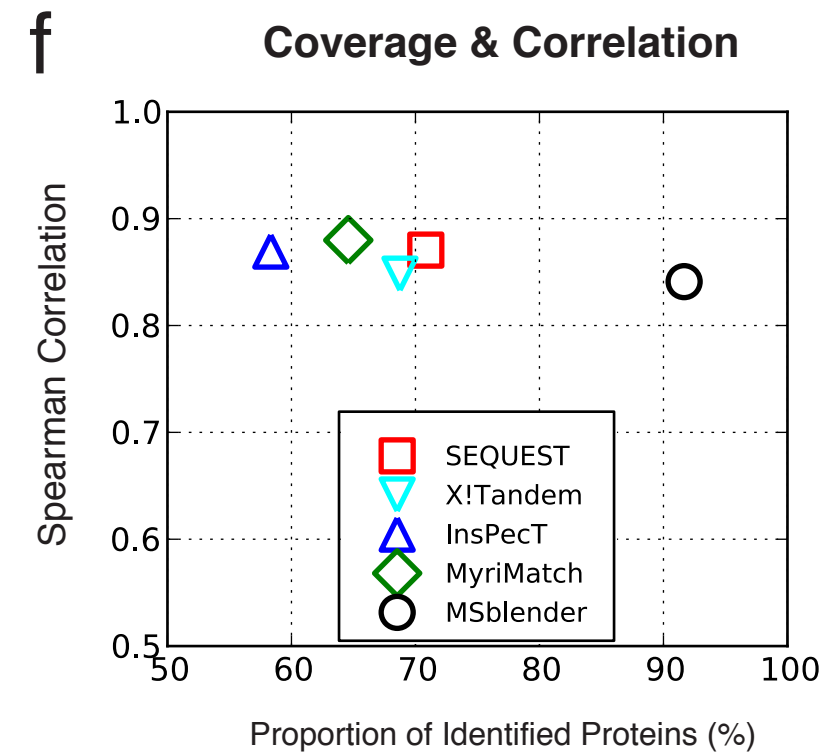
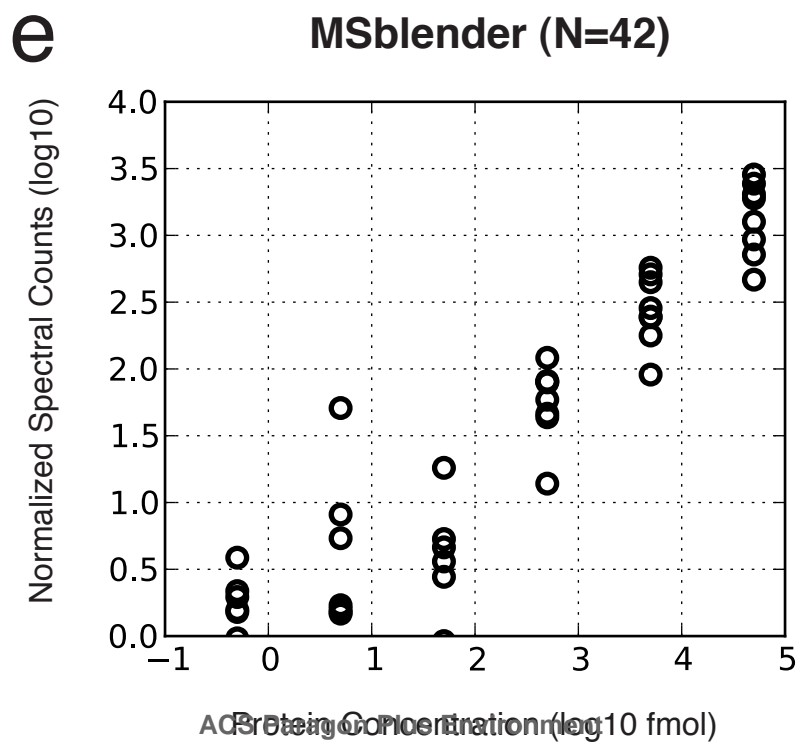
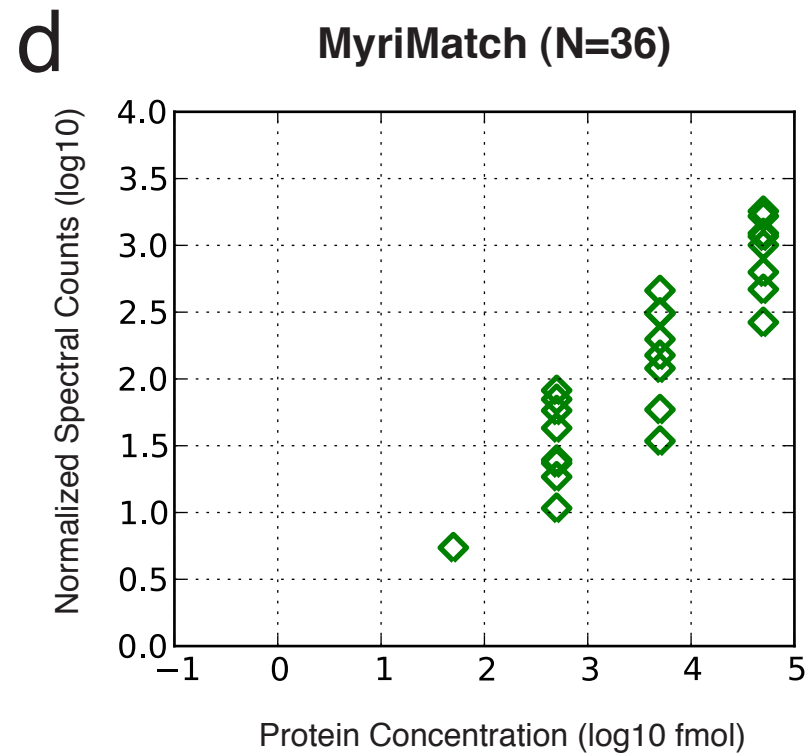
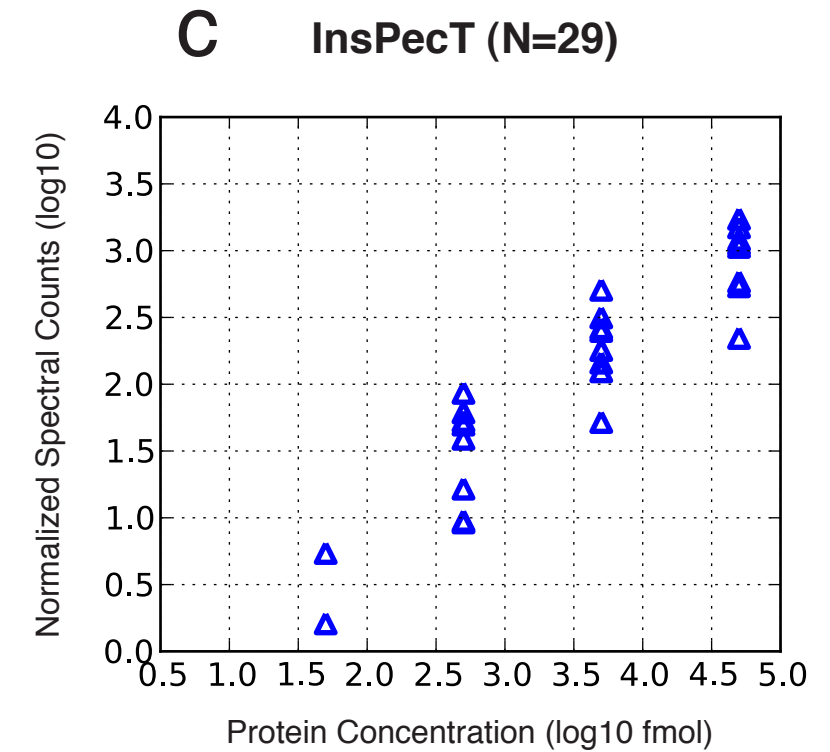
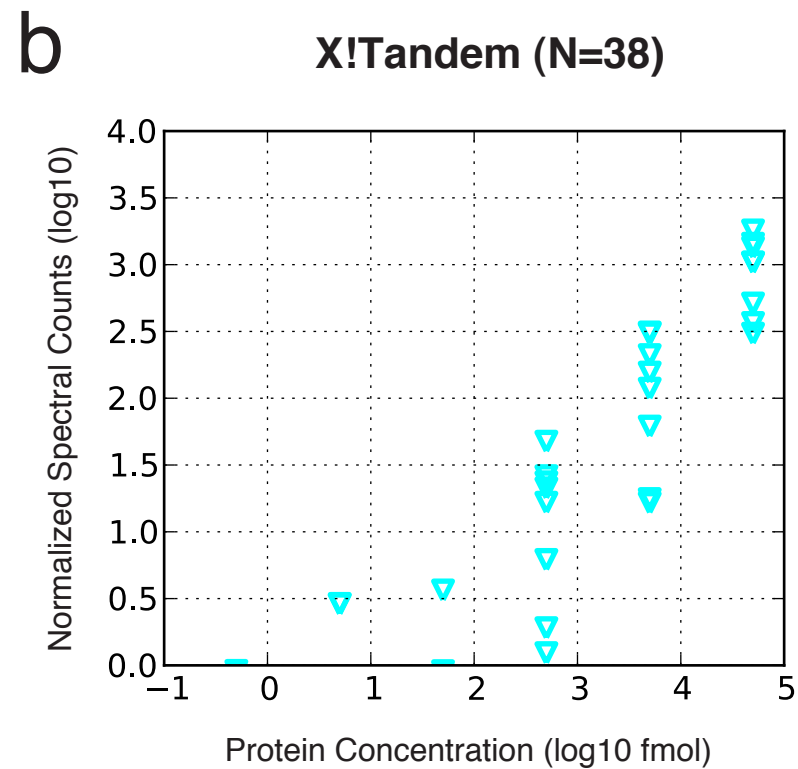
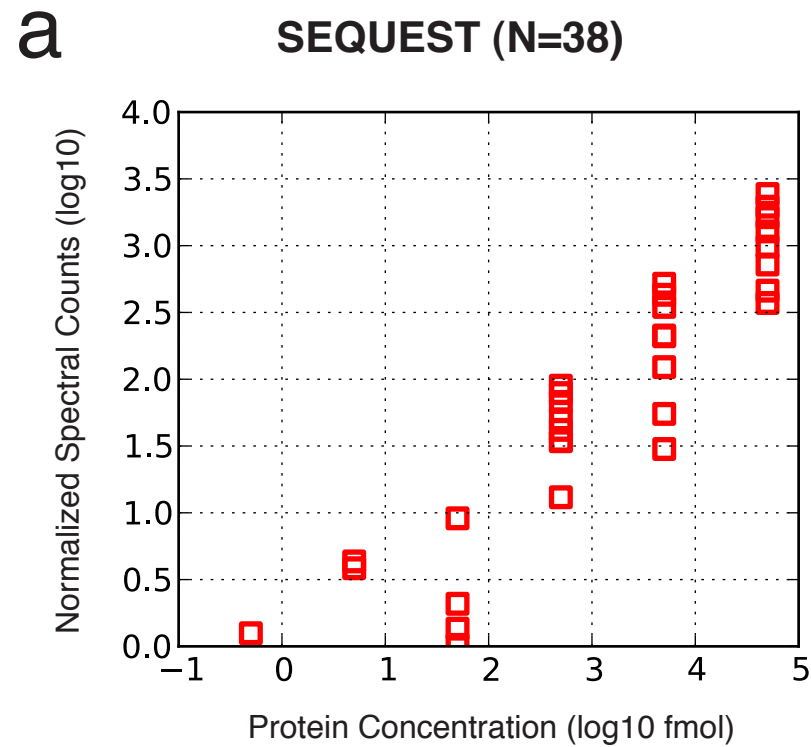




## Accuracy of FDR

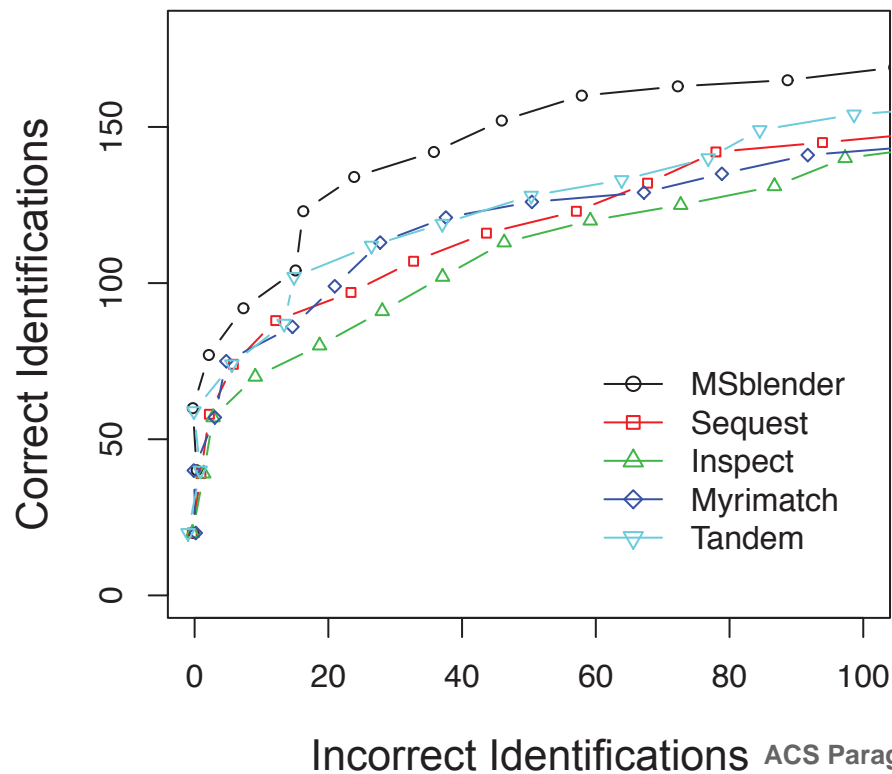
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24





**a**

## Raw Spectral Counts

**b**

## Split Spectral Counts

