

New Method for Label-Free Quantification in the Proteome Discoverer Framework

David M Horn,¹ Torsten Ueckert,² Kai Fritzscheier,² Katja Tham,² Carmen Paschke,² Frank Berg,² Hans Pfaff,² Xiaoyue Jiang,¹ Shijun Li,¹ and Dani Lopez-Ferrer,¹
¹Thermo Fisher Scientific, San Jose, CA, USA; ²Thermo Fisher Scientific, Bremen, Germany

ABSTRACT

A new label-free quantification method based on the Minora algorithm is presented and compared to pre-existing label free quantification methods in the Thermo Scientific™ Proteome Discoverer™ software framework. The results of the new algorithm were significantly more accurate across a wide dynamic range compared to spectral counting and "Top N" quantification. The new algorithm was also run on a subset of the Akhilesh Pandey human proteome dataset to identify proteins specific to specific tissue types.

INTRODUCTION

Proteome Discoverer software is a node-based workflow engine and study management platform for analysis of mass spectrometry-based proteomics datasets. The latest released version 2.1 fully supports isotopically-labeled quantitative workflows, such as TMT™ reporter ion-based quantification and SILAC precursor ion quantification, but the supported label-free quantification methods are significantly less sophisticated. Currently, spectral counting is possible but not recommended when quantitative accuracy is required. The only supported label-free quantification workflow produces an average abundance of the top "N" most abundant peptides and this has been shown to be accurate for even highly complex datasets. However, "Top N" quantification results cannot be used to create ratios, scaled abundance values, or to be used as replicates to generate standard errors. Here we present a new workflow for untargeted label-free quantification using a new feature detection approach that provides the full suite of quantitative capabilities previously only available for isotopically-labeled quantification. The workflow will be compared to the two aforementioned label-free quantification workflows available within Proteome Discoverer 2.1 software.

MATERIALS AND METHODS

A standard dataset of Arabidopsis proteasome proteins spiked into a background of *E. coli* proteins (PXD003002) was downloaded from the PRoteomics IDentifications (PRIDE) repository. This dataset was originally used to evaluate a spectral counting algorithm and is described in reference 1. The Pandey human proteome dataset² was also downloaded from PRIDE and a portions of the dataset to demonstrate untargeted label free quantification of data with a multi-dimensional separation.

For quantification using spectral counts, each of the datasets with the different levels of *Arabidopsis* proteasome proteins was run separately in batch mode using a standard Sequest™ HT-Percolator workflow and a basic consensus workflow. Subsequently, all Processing results were reprocessed using a single Consensus workflow with the "Merge Mode" parameter in the MSF files node set to Do Not Merge. With this setting, the number of unique peptides and PSMs for each of the datasets will be represented as a separate column. The Sequest HT search was performed against the entire *Arabidopsis thaliana* and *Escherichia coli* databases. The table with PSM values for each sample was exported to Microsoft Excel format and ratios were calculated manually.

For the "Top N" quantification workflow, a Precursor Ion Area Detector node was incorporated in the Processing workflow used for spectral counting above. The default "CWF_Comprehensive_Enhanced_Annotation_Quan" template was used for the Consensus workflow. In the Peptide and Protein Quantifier node, the "Top N Peptides Used for Quantification" parameter was set to 3. Like for spectral counting, the table with the reported Top N protein abundances was exported to Excel and ratios were calculated manually.

New Method for Feature Detection

The new feature detection algorithm is an extension of the Minora algorithm, which had already been used for precursor ion quantification since the release of Proteome Discoverer 1.2 software. Minora had always detected all isotopic peaks in a given data set, but up to now only those LC/MS peaks associated with peptide spectral matches (PSMs), and their associated isotopic forms in the case of SILAC, were used for quantification. In this pre-release version of Proteome Discoverer 2.2 software, the Minora algorithm has been modified to detect and quantify isotopic clusters regardless of whether or not they are associated with a PSM.

A typical Processing workflow for Minora feature detection is shown in Figure 1. The new label-free quantification workflow can be invoked by simply attaching the "Minora Feature Detector" to the Spectrum Files node. This new feature detector will also be used for the isotopically-labeled precursor quantification method such as SILAC.

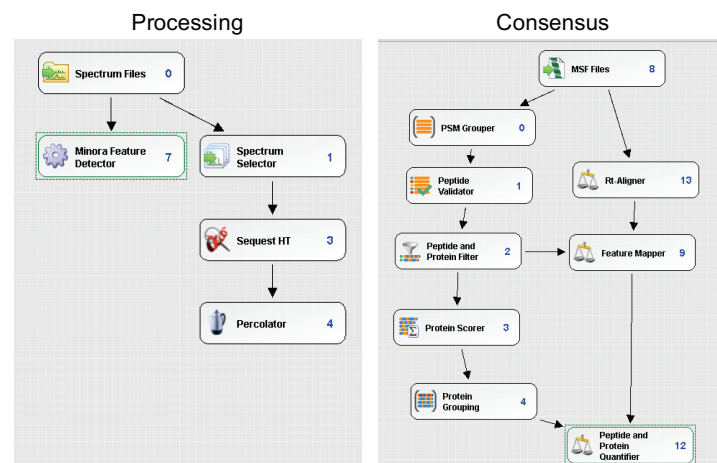


Figure 1. Typical Processing and Consensus workflows for untargeted label-free quantification. The Minora Feature Detector, Rt-Aligner and Feature Mapper are new nodes created for the untargeted label-free quantification workflow. The Minora Feature Detector will also replace the old Precursor Ions Quantifier node used for SILAC and other precursor ion quantification workflows.

A typical Consensus workflow for label free quantification is also shown in **Figure 1**. There are two new nodes added to this workflow that perform retention time alignment and feature mapping. The feature mapper groups features detected from the Processing runs into "Consensus Features" that are mapped and quantified across all raw files and performs gap filling to find features that were not initially detected in the processing workflows. The Peptide and Protein Quantifier node works as previously, with improvements to scaling and normalization that benefit all quantification workflows.

There are three new tabs for feature detection results in the consensus report: Consensus Features, LCMS Features, and LCMS Peaks. The LCMS features are isotopic clusters grouped together for a given raw dataset and consist of multiple LCMS Peaks. Ultimately, the release may not include the LCMS Peaks list given that as much as 10's of millions peaks could be detected in complex datasets. The consensus features link directly to the associated peptide group as well as the list of LCMS features detected from each data files (**Figure 2**). Also, when a consensus feature is selected, the traces for each of the features are shown in the chromatogram traces view. When a single LCMS Feature is selected, the chromatographic profile for only that individual feature is displayed.

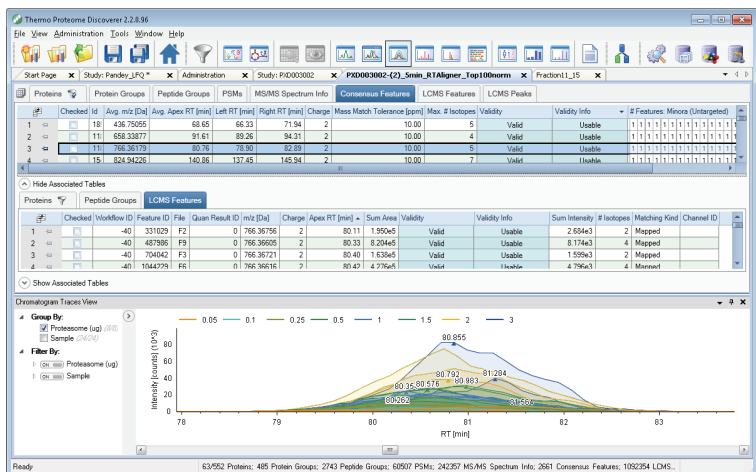


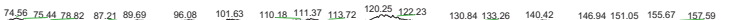
Figure 2. The Consensus Features table is linked to the collection of LCMS Features from each raw file. The chromatographic profiles for each LCMS Feature are shown in the Chromatogram Traces View at the bottom.

Like the other quantification workflows in Proteome Discoverer 2.1 software, the peptide group abundances from the new label-free quantification method are calculated as the sum of the abundances of the individual PSMs for a given study factor that pass a quality threshold. The protein abundance is calculated as the sum of the peptide group abundances associated with that protein.

RESULTS AND DISCUSSION

The database searches produced a total of 55 *Arabidopsis* proteins and 423 *E. coli* proteins. This is less than would be expected given the relatively high protein concentration and long gradient length, but the chromatography used for these data analyses was suboptimal with peaks up to 5 minutes wide (**Figure 3**). Also, as the amount of the *Arabidopsis* proteins added to the sample increased past 1 µg, the peptides from these proteasome proteins dominate the chromatogram and thus the number of *E. coli* proteins decreases dramatically with increasing *Arabidopsis* protein concentration (data not shown).

0.05 µg proteasomes
1 µg *E. coli*



0.5 µg proteasomes
1 µg *E. coli*



3 µg proteasomes
1 µg *E. coli*

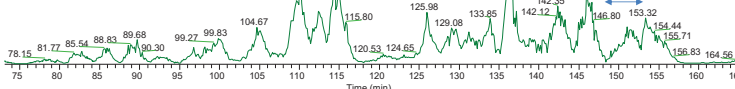


Figure 3. Base peak chromatograms for three of the LC/MS runs, each scaled to 2e7 intensity. The dataset at the bottom is dominated by *Arabidopsis* peptides, leading to significant suppression of the *E. coli* peptides. Also, the typical chromatographic peak in this chromatogram can be up to 5 minutes wide, also decreasing the number of peptides and proteins that can be identified.

The spectral counts-based quantification results correctly indicate the direction of expression for the *Arabidopsis* proteasome proteins, but the ratios are inaccurate for the more extreme ratios. The response is also relatively non-linear, with the average ratio for the 0.1 µg/1 µg samples showing a lower value than the 0.05 µg/1 µg samples and the 3 µg/1 µg ratio measuring lower than the 2 µg/1 µg ratio. These results are not a surprise given that it is widely known that this type of spectral counting is not expected to produce accurate quantification results. Normalized spectral counting algorithms are a significant improvement over the basic spectral counting method shown here and reference 1 from which these data were obtained describes a such a method. Implementation of such a method using emPAF values is planned for the individual study factors is being considered for a future Proteome Discoverer software release. However, all spectral counting-based quantification methods usually provide poorer sensitivity and dynamic range than other quantitative techniques due to the requirement for multiple PSMs for any given protein. As can be seen in Table 1, less than half of the *Arabidopsis* proteins could be quantified across the full dynamic range due to lack of PSMs in the samples with lowest protein abundance.

The "Top N" protein quantification results are shown in the second set of columns in Table 1. The accuracy of the ratios is noticeably improved compared to spectral counting, producing a response that is closer to linear. However, there are fewer quantified proteins in the "Top N" method than for spectral counting, primarily due to the requirement that the same three peptides need to be identified across all of the datasets. This is in effect even more stringent than the spectral counting method above and as a result even fewer proteins are quantified across the samples. Also, while the accuracy of the ratios is improved, the precision of the measurements are not much improved over spectral counting.

For feature detection-based quantification, the calculated ratios were significantly closer to the theoretically expected values at the lowest *Arabidopsis* concentrations. The precision of the ratios was also significantly improved in almost all cases for the feature detection results. The use of feature mapping led to a significantly increased number of quantified proteins given that only a single PSM is required for a given peptide across all raw files. The accuracy and precision of this method also benefits from the use normalization based on the *E. coli* proteins, which are known to be equally abundant across all samples.

A screen shot of the *Arabidopsis* protein identification results with untargeted label-free quantification is shown in **Figure 4**. The ratios and the scaled abundances for the identified proteins and peptide groups are color coded based on the level of expression. Scaled abundances were originally introduced in Proteome Discoverer 2.1 software primarily for the TMT quantification workflow and are now available in the preliminary version of Proteome Discoverer 2.2 software for feature detection-based label free quantification. The samples can be sorted by scaled abundance for any given sample type, as seen in **Figure 4** for the highlighted 0.05 sample group. It can be easily seen that each of the proteasome proteins exhibit a similar trend by simply looking at the pattern of blue and red boxes. Also, since the scaled abundances exhibit the same profile as the ratios, the need for the calculation of ratios is somewhat obviated.

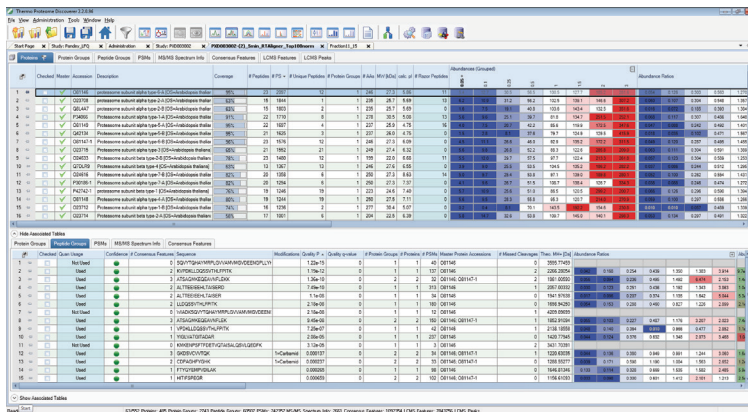


Figure 4. Untargeted label-free quantification results within the Proteome Discoverer software framework. Both the ratios and the scaled abundance values are color-coded to display significantly under- or over-expressed proteins.

