# MassHunter Profinder:
# Batch Processing Software for High Quality Feature Extraction of Mass Spectrometry Data

## Technical Overview

## Introduction

LC/MS metabolomics workflows typically involve the following steps: (1) data acquisition using high resolution accurate mass LC/MS, (2) feature extraction from raw data, (3) statistical analysis, (4) annotation and identification using database and library matching, and (5) increasingly, the use of pathway analysis.

The feature extraction step plays a critical role in differential analysis, as the quality of the feature extraction results impacts subsequent steps in the metabolomics workflow. Ultimately, it affects the effectiveness with which raw data is transformed into biologically relevant information.

Previously, a two-pass feature extraction[1] process was used to find compounds from complex accurate mass data. Compounds are comprised of neutral mass, retention time (RT), abundance and $m/z$. The first pass used Molecular Feature Extraction (MFE) in MassHunter Qualitative Analysis (Qual). Next the results were imported into Mass Profiler Professional (MPP) for binning, aligning, and creating a consensus for each compound. The second pass used the list of consensus compounds and the "Find by Ion (FbI)" algorithm in Qual for targeted feature extraction. This two-pass feature extraction approach improved the quality and accuracy of mass, retention time and abundance values for each candidate compound compared with the MFE only data mining. However, this workflow has some limitations such as low throughput, use of two separate software programs, as well as a lack of data visualization and curation tools.

MassHunter Profinder is a new software tool designed for batch processing of large, complex accurate mass LC/MS data. Profinder also integrates many software functions into one dedicated processing tool. As a result, this software reduces the complexity of raw data and improves the quality of extracted compounds. Most importantly, it improves user experience and increases the throughput using fully automated batch data processing.

**Agilent Technologies**

This Technical Overview describes the overall data processing workflow in Profinder, including a detailed review of the recursive algorithms, the simplified GUI design, as well as the benefits of grouping and batch analysis, which was not possible before. It also demonstrates the utility of Profinder using a relatively large batch of yeast metabolomics data that were acquired in a nontargeted LC/Q-TOF MS approach[2]. The improved results obtained by Profinder are illustrated through quality enhancement using statistical analyses.

## Key Functionality and Benefits

- Supports untargeted and targeted data mining of LC/TOF and LC/Q-TOF accurate mass data

- Integrates MFE and FbI in a single workflow, eliminating multiple import/export of results between Qual and MPP

- Increases data processing throughput using batch feature extraction

- Results in greater recovery of missing features by reducing false negatives through recursive feature extraction algorithm

- Offers compound-centric visualization across multiple sample files, enabling quick sorting and filtering of compound group results, visual review and identification of outliers for manual editing

- Extracted ion chromatograms and mass spectra are optionally overlaid and colored by sample group, enhancing visual inspection and comparison of two or more sample groups with or without replicates

# Batch Feature Extraction Workflows

There are three feature extraction workflows.

**Batch Recursive Feature Extraction** is the primary untargeted feature extraction workflow in MassHunter Profinder. This workflow integrates the batch MFE and batch FbI two-pass feature extraction workflow into a single processing tool, eliminating the extra effort required before. A wizard guides the user through parameter setup. The data input and parameters settings take only a few minutes, and now the two-pass feature extraction process is automated and self-contained in Profinder.

Figure 1 shows the results of Profinder Batch Recursive Feature Extraction for yeast metabolomics data acquired on an Agilent Q-TOF LC/MS in positive ion mode. The main view consists of four linked navigation windows. The Compound Group Table (Figure 1A) displays the compound-level information grouped and summarized across multiple data files. The individual file details of a selected compound group are shown in the Compound Details Table (Figure 1B), Extracted Ion Chromatogram (Figure 1C), and MS Spectrum (Figure 1D) windows. This compound centric visualization allows a detailed inspection of feature extraction results. Users can quickly sort compound groups in multiple ways (i.e. neutral mass, RT, abundance, found, and missing) to identify spurious features for deletion. In the Chromatogram and MS Spectrum windows, the plots can be colored by individual sample file (Figure 2A), or by sample group (Figure 2B and 2C) for enhanced visualization. Furthermore, EICs can be stacked by individual sample file (Figure 2A and 2B), overlaid by sample groups (Figure 2C), or all the files can be overlaid (Figure 2D). This function is useful for quick comparison of a compound feature across sample files or sample groups, visualizing missing peaks, and manually editing EIC peak integration.
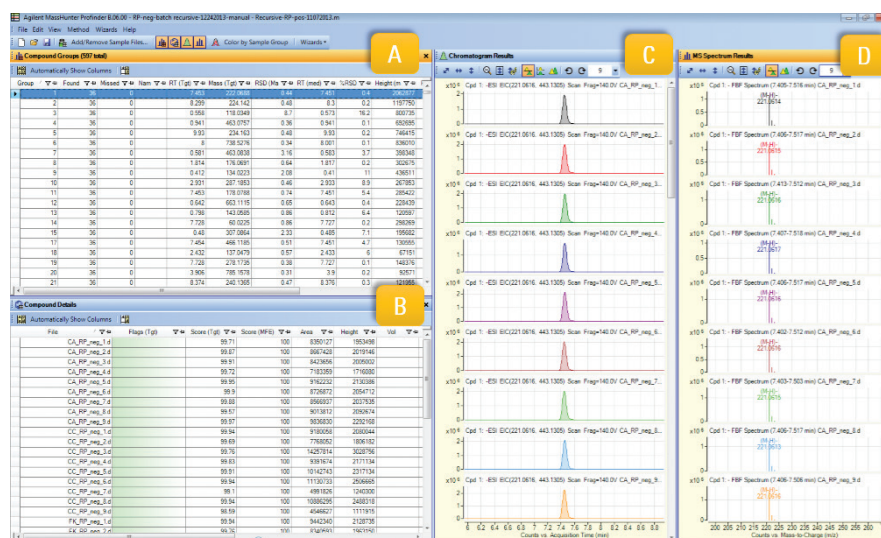


Figure 1. MassHunter Profinder main view shows the results using Batch Recursive Feature Extraction. There are four windows: (A) Compound Group Table, (B) Compound Details Table, (C) Extracted Ion Chromatograms (EICs), and (D) Mass Spectra.
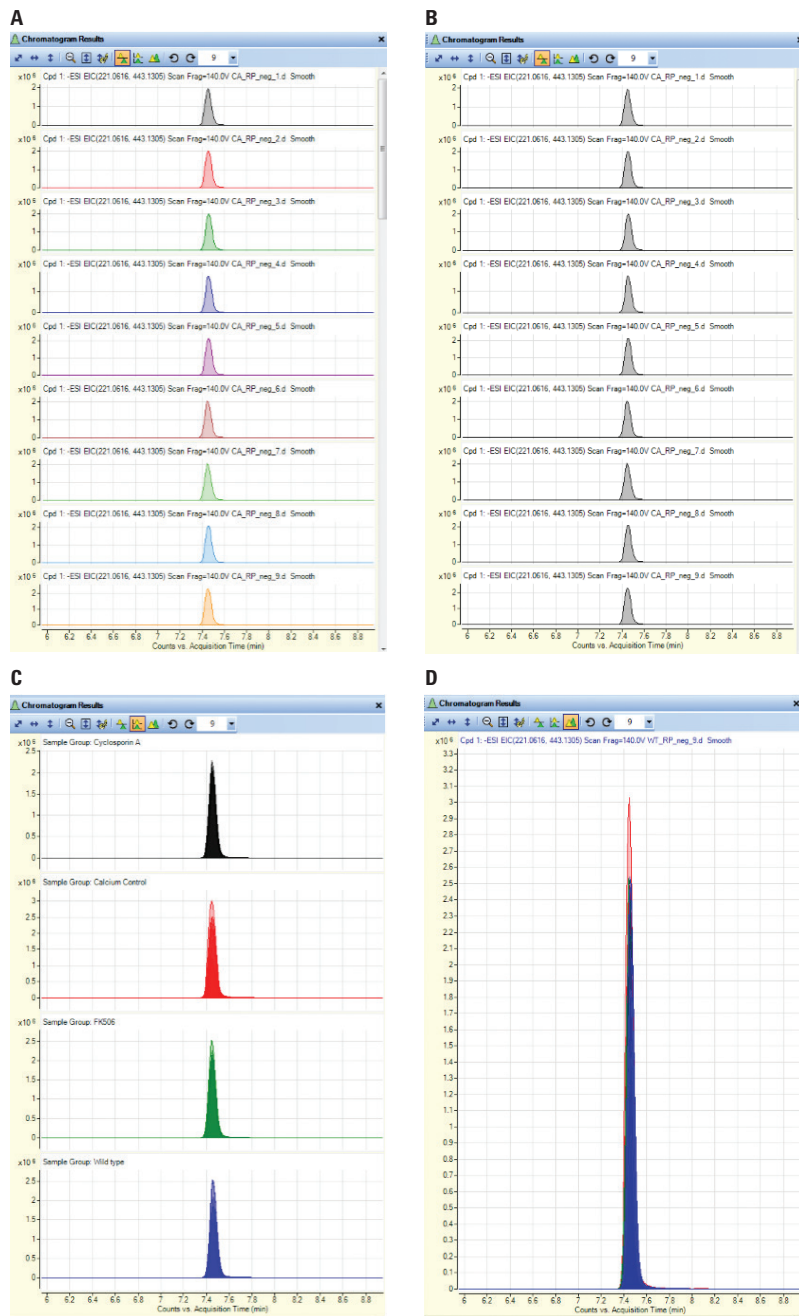
Figure 2. EICs visualization modes (A) List mode, color by files, (B) List mode, color by sample group, (C) Sample group mode, color by sample group, and (D) Overlay mode.

Batch Recursive Feature Extraction provides greater missing feature recovery. This is achieved by binning and alignment of compound features in the first-pass MFE, after which a composite spectrum for all found ions is created for each consensus (summary) feature. The composite compound feature list is then used as a target list for the second-pass FbI feature extraction.

Batch Recursive Feature Extraction will sometimes report a zero abundance value for a compound. Figure 3 illustrates how a selected compound in one of four sample groups was curated. The EICs (Figure 3A, red traces) were extracted, however, their peak abundances were below the user-specified filter criteria and the peak abundances were reported as zero. Profinder provides raw data visualization and a manual integration tool, a very important feature of this software, which enables visual validation of the feature extraction results and manual editing of EIC peaks. As a result, the zero abundance values of this compound group were manually corrected for this sample group (Figure 3B, red traces).
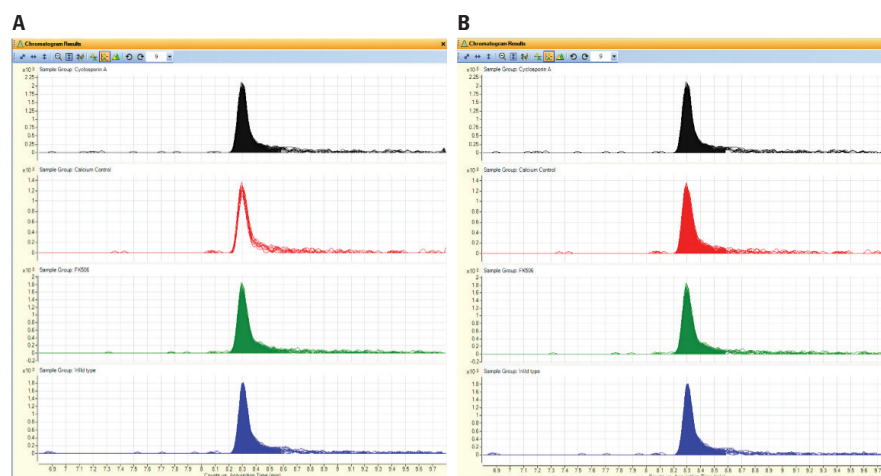
**A**  **B**



Figure 3. Correcting a missing compound. EICs are overlaid by sample groups, before (A) and after (B) manual peak integration.

**Batch MFE**[3] uses a new proprietary algorithm to extract compounds from a batch of raw data files. Compounds from multiple files are aligned and binned using neutral mass and retention time. As a part of the alignment step, the software automatically reviews and regroups assigned compound ions using batch information. This workflow is highly beneficial for optimizing the settings for the primary Batch Recursive Feature Extraction workflow. This is because Batch MFE is a high speed untargeted feature extraction algorithm and is faster than Batch Recursive Feature Extraction.

**Batch Targeted Feature Extraction** allows the user to extract compounds of interest with known chemical formulas from large complex data sets. As shown in Figure 4, the user has several options, one can input formulas directly, supply a list of formulas by using a compound exchange file (.CEF), or use a chemical compound database (.CSV or Agilent .CDB format). This workflow offers advantages of high selectivity, fast data processing, and tentative compound annotation. It also provides a useful tool for biological pathway-driven data analysis through Targeted Feature Extraction using a database, for example, generated from Agilent MassHunter Pathways to PCDL software. The Pathways to PCDL software is designed to create an organism specific metabolite database by selecting pathways of interest from public sources BioCyc, KEGG, or WikiPathways. This software was used to create a folic acid biosynthesis pathway database based on the WikiPathways data source for *Saccharomyces cerevisiae*. Using the folic acid biosynthesis pathway database and Batch Targeted Feature Extraction workflow, five compounds (glutamate, ADP, GTP, L-serine, and phosphate) (Figure 5) were extracted and annotated from the yeast metabolomics data acquired on an Agilent Q-TOF LC/MS in negative ion mode.
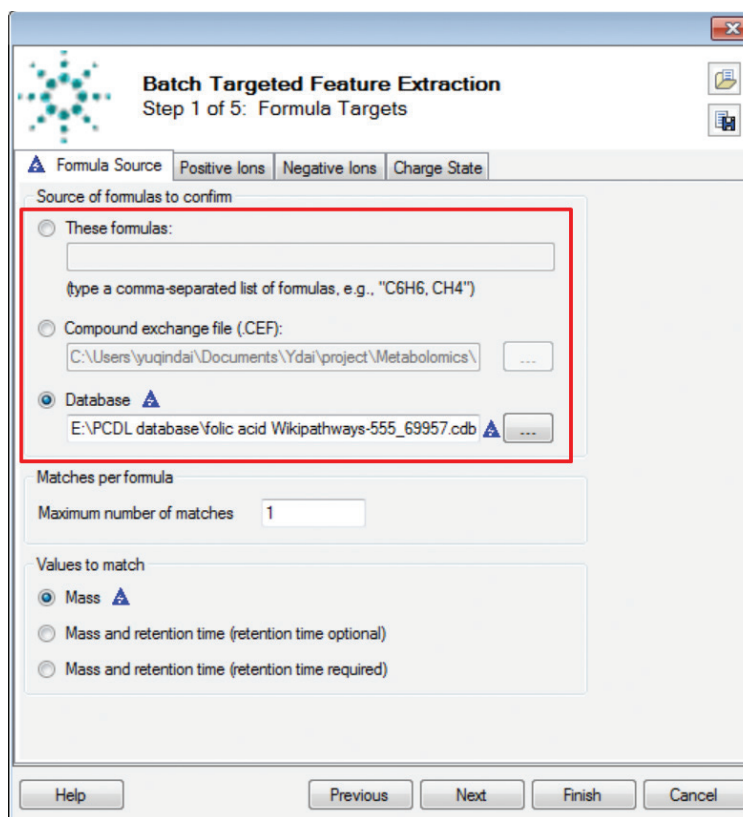


Figure 4. Three different ways to supply formulas in Batch Targeted Feature Extraction.

Figure 5. Batch Targeted Feature Extraction results from the yeast metabolomics data.

Consistent peak integration across all files or sample groups is important for reliable downstream statistical analyses. When isobaric compounds in a complex matrix elute closely, it is sometimes difficult for peak integrators to establish the peak start/end. As illustrated in Figure 6, the EIC peaks were not consistently integrated across all data files (Figure 6A). Profinder allows the user to investigate EICs and perform manual integration consistently across all overlaid files. The curated peak abundances provide more reliable information for comparison between sample groups or sample files (Figure 6B).

**A**                                          **B**



Integration across all files

Figure 6. Reviewing for consistent peak integration across all overlaid files, before (A) and after (B) manual correction of peak integration.

**Profinder Batch Recursive Feature Extraction Improves Quality of Statistics Analyses**

Profinder has a significant impact on statistics from mass spectrometry data. By implementing two-pass batch feature extraction, data visualization, and a manual curation tool, Batch Recursive Feature Extraction reduces the number of false positives and false negatives, subsequently decreasing CVs within sample group replicates. To demonstrate this, we compared the results from Batch Recursive Feature Extraction with Batch MFE in MPP (Figures 7–9)

Figure 7 shows two histograms of CV % within compound groups (that is, the error in measured abundances for compounds aligned by mass and RT across multiple sample files) using Batch Recursive Feature Extraction (Figure 7A) and Batch MFE (Figure 7B). There were 36 samples across 4 sample groups. The lower CV % values suggest that the abundance measurements within a given compound group are of higher quality. The higher CV % values suggest that missing or incorrect peak integrations may be present within a given compound group. Two observations can be noted from the plots. First, the histograms are clearly shifted to lower CV % using Batch Recursive Feature Extraction as compared to Batch MFE. Second, the cumulative percentage of compound groups with low CV % values increased as the histogram distribution shifted to lower variance rates. These trends suggest an overall improvement in compound quality within those respective compound groups.
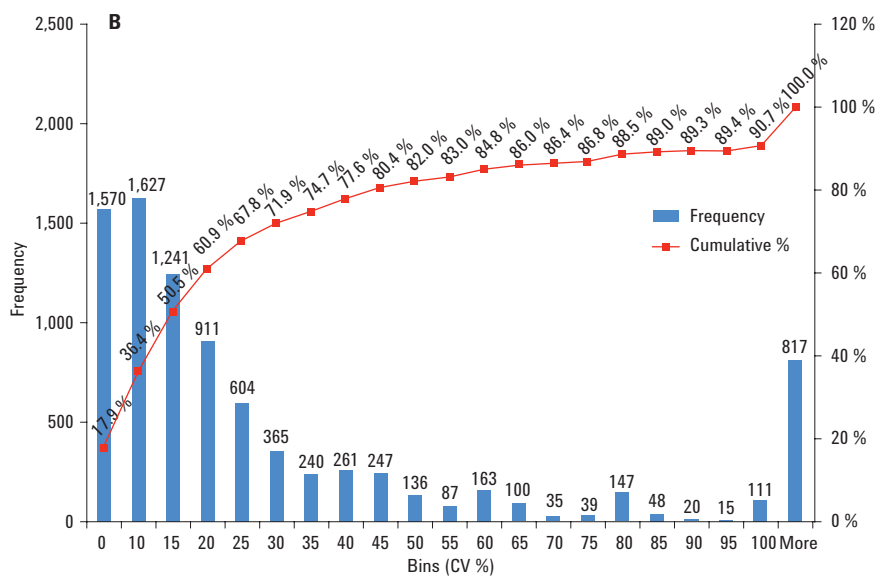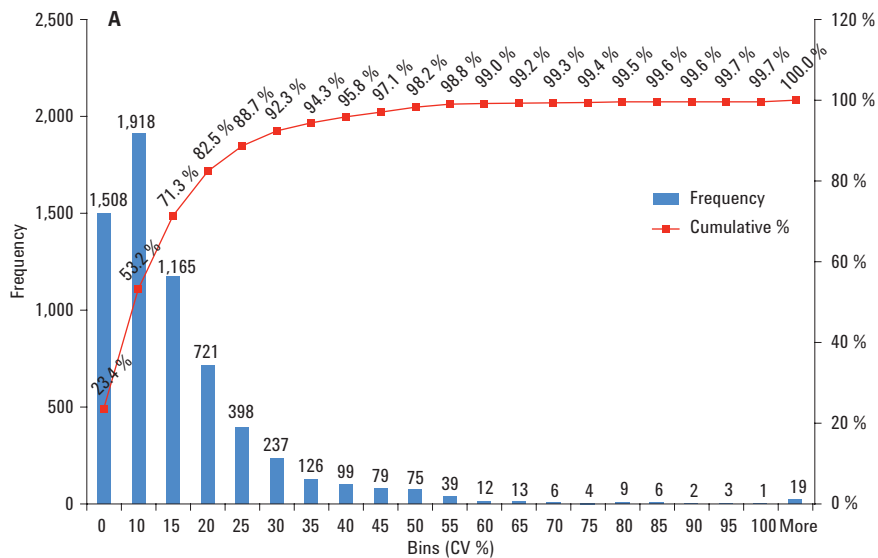
Figure 7. Distributions of binned compound CVs and their frequency (A) Batch Recursive Feature Extraction and (B) Batch MFE. The red line represents the cumulative percent of compounds at a given CV.

Figure 8 shows frequency histograms for 36 sample files extracted by Batch Recursive Feature Extraction (Figure 8A) and Batch MFE (Figure 8B). "Frequency" in MPP refers to the number of files in which a compound was found within a given compound group. The maximum frequency is, therefore, 36, given 36 sample files. If there are missing compounds as a result of the feature extractor not finding a compound within a given file, the frequency will be less than 36. Following Batch Recursive Feature Extraction, the compound groups have a full 36/36 compound representation (Figure 8A).
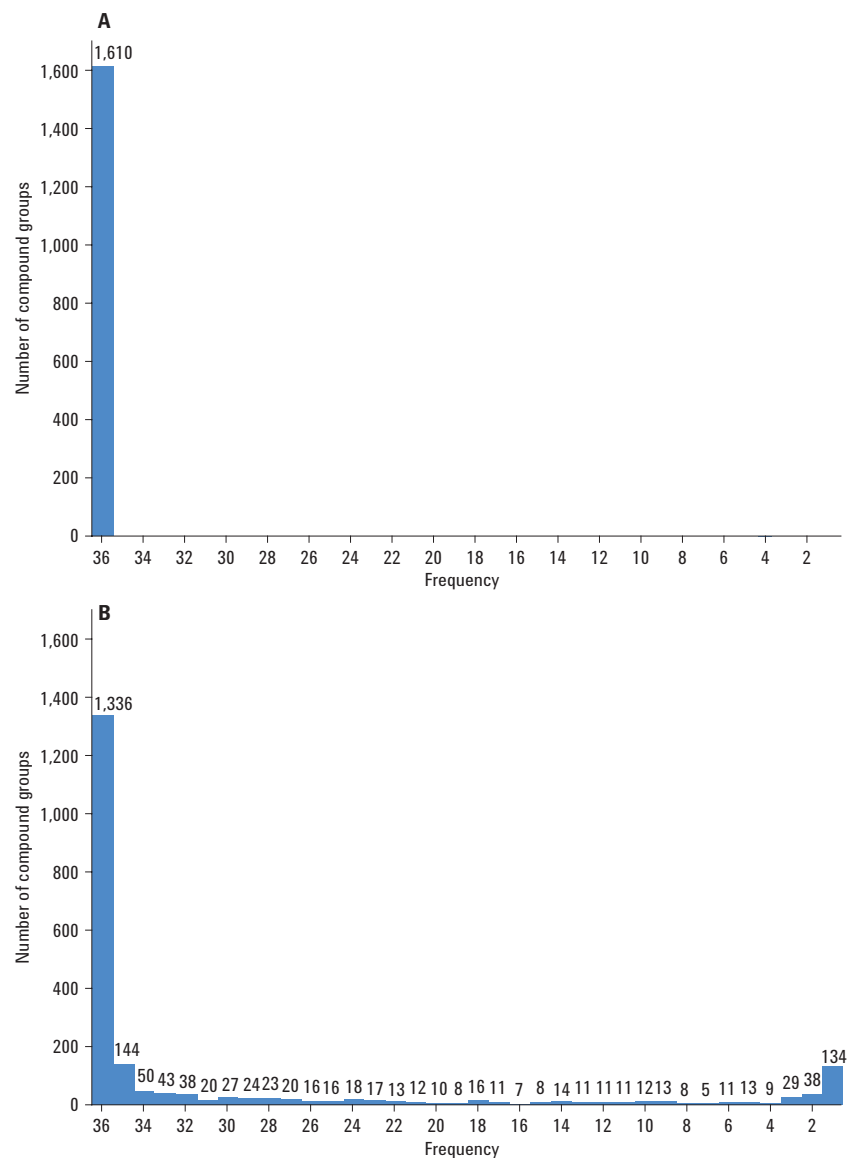


Figure 8. Frequency distribution of extracted compounds across all samples (A) Batch Recursive Feature Extraction and (B) Batch MFE.

10

Figure 9 shows 2D Principal Components Analysis (PCA) plots colored by sample group for 36 individual sample files representing four sample groups. It is clear that Batch Recursive Feature Extraction has improved the separation between sample groups and also tightened the clustering of the sample replicates. This suggests a marked reduction in noise within sample group replicates. In turn, there is greater visibility of the true sample group-dependent variance.

It is evident that, compared to the results from the Batch MFE only, Profinder Batch Recursive Feature Extraction significantly improves the quality of all subsequent statistical analyses, as evidenced by the higher frequency of lower CV % compound groups in the histograms (Figure 7), complete representation of compounds in each compound group (Figure 8), and better separation between the four biological groups in PCA plots (Figure 9).
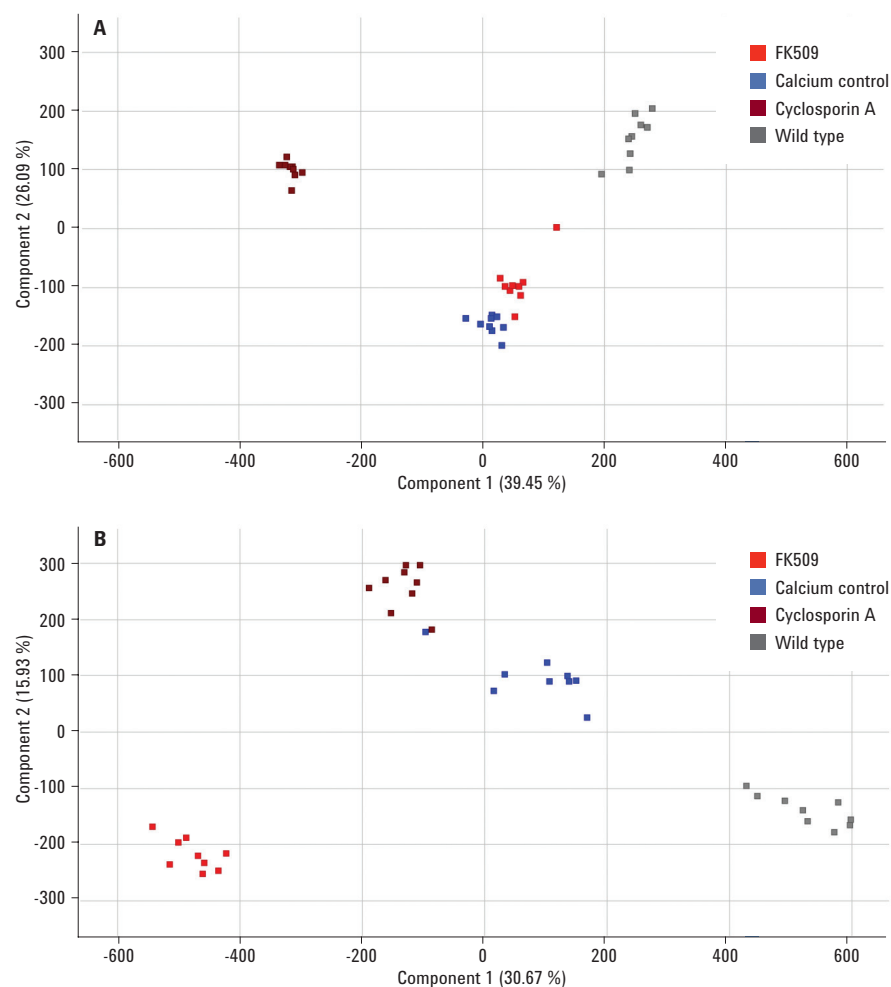


Figure 9. 2D PCA plots of 36 samples in four groups (A) Batch Recursive Feature Extraction and (B) Batch MFE.

## Conclusions

MassHunter Profinder is a powerful feature finding software for high resolution accurate mass LC/MS data. It enables fully automated and self-contained two-pass batch feature extraction from large complex data sets. The intuitive and easy-to-use MassHunter Profinder software significantly improves the quality of data for reliable statistics analyses. This software can also be used to support applications beyond metabolomics, such as proteomics and food profiling.

## References

1. N. Kitagawa, *et al.* Improving Untargeted Differential Analysis of Mass Spectrometric Data by Recursive Feature Extraction. ASMS (**2009**).

2. S. Jenkins, *et al.* "Compound Identification, Profiling and Pathway Analysis of the Yeast Metabolome in Mass Profiler Professional" Agilent Application Note, publication number 5991-2470EN (**2013**).

3. N. Kitagawa, *et al.* A Novel Two-pass Feature Extraction Workflow for the Statistical Profiling of Mass Spectrometric Data. ASMS (**2013**).

www.agilent.com/chem

**Agilent Technologies**