**Agilent**

Trusted Answers

# Processing Agilent AutoMSMS Proteomic Data Using Mascot Online

## Introduction

One of the key challenges of proteomics is the fast and accurate identification of proteins. This process uses spectral matching of AutoMSMS data to known protein sequences, together with an Agilent LC/Q-TOF system using Agilent MassHunter software and Mascot MS/MS Ions Search (Matrix Science). These tools enable the creation of a proteomic workflow.

Acquiring good quality data is important. Data should be checked for chromatography, mass accuracy, fragmentation, and confounding artifacts such as excessive sodium adducts or undigested protein. This can be done using Agilent MassHunter Qualitative Analysis software. The trypsin digestion protocol should be checked for robustness and reproducibility. Data should be acquired as centroid to transition smoothly to the data extraction step.

Data extraction produces a reduced data set in the form of an MS/MS peak list text file with additional information, including precursor $m/z$ and retention time. The aim is to export signal and reduce noise. The size of the extracted file will affect the proteomic search time. MassHunter Qualitative Analysis software version 10.0 can be used. Vendor-neutral data extraction software is available, such as MSConvert from ProteoWizard. Care should be taken to optimize data extraction parameters using an iterative approach.

The database search will first try to match tryptic peptide masses to a protein sequence database, including any modifications selected. It will then try to match the acquired extracted MS/MS data to a theoretical fragmentation of the matched peptide sequences and assign a score. This is the search space, which puts a high demand on computational resources. Care should be taken to optimize parameters to reduce search time without impacting data quality. This can be done by optimizing the size of the data file for signal-to-noise and reducing the search space by using both a smaller database and a smaller MS ppm error window.

Results should be checked for correct modifications, including: cysteine with and without modification, number of unexpected missed tryptic cleavages, mass accuracy, score, and percent sequence coverage. The use of a system check sample (similar to your target) is highly recommended. Reporting can be as simple as copying and pasting, or exporting the results to a .CSV file.

# Workflow overview of Mascot (online)

Mascot online from Matrix Science is a free version of the Mascot Server 3.0 and related products (note that the free version is limited in terms of function and upload file size). It is frequently used as an introduction to proteomic data processing due to its easy-to-use interface and online tutorials. Agilent AutoMSMS proteomic data can be searched using the Mascot MS/MS Ions Search page.

For data extraction using Agilent MassHunter Qualitative Analysis 10.0 through a mascot generic format file (.MGF), choose the following:

1. Select **Configuration** and click **Show Advanced Settings**.

2. Start from a Default.m method and save as DA_Qual_MFG_Export.m.

3. Using the **Compounds** view, navigate to **Method Editor Export MGF options**. Extract the entire data file (Figure 1). This tends to give the highest Mascot score. If the file is too large, it can be preprocessed with **Method Editor Compound Discovery Find By Auto MS/MS** using the default parameters (make sure that **Results Extract Complete Result Set Automatically** is unticked and that **Highlight First Compound** is selected to reduce processing time). If the file is still too large, the **Limit to the largest 10,000 Compounds** function can be used. If a preprocess is used, then select the **Export the Extracted Spectra** function (Figure 1).
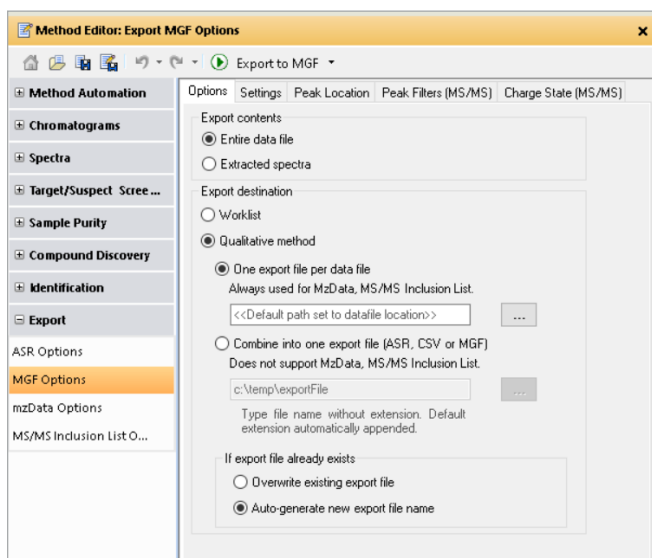


**Figure 1.** Method Editor: Export MGF Options > MGF Options > Options.

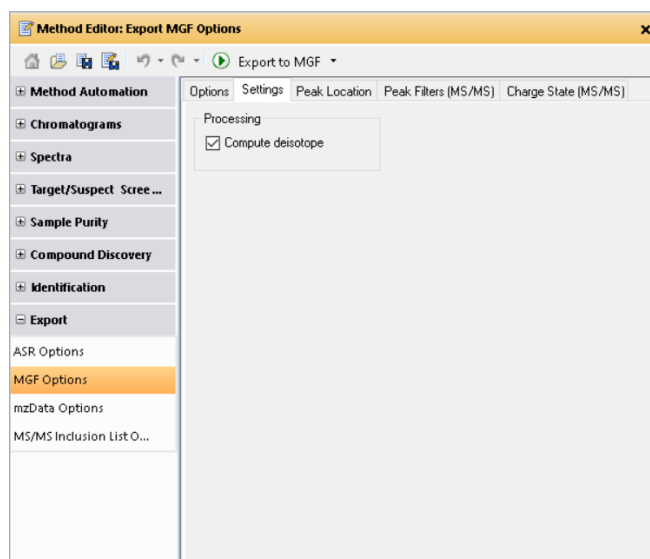4. Under **Method Editor Export MGF options**, select **Compute deisotope** (Figure 2).



**Figure 2.** Method Editor: Export MGF Options > MGF Options > Settings.

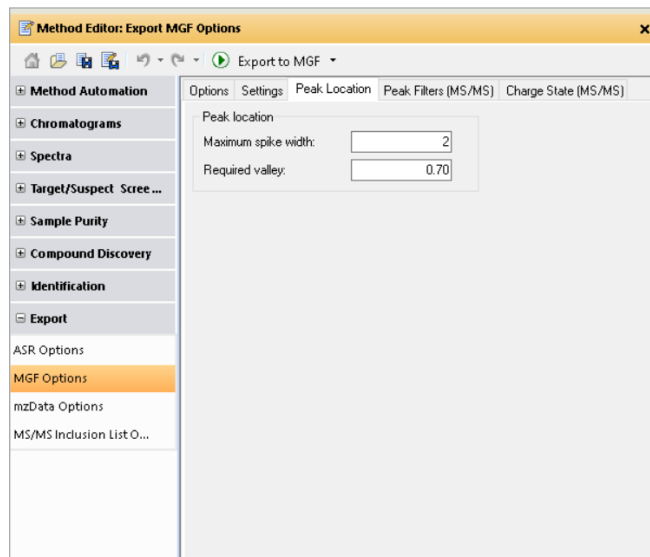5. For **Peak Location**, leave as default (Figure 3).



**Figure 3.** Method Editor: Export MGF Options > MGF Options > Peak Location.

6. For **Peak Filters (MS/MS)**, deselect both height filters, and set Maximum number of peaks to 50 to export signal and reduce noise (Figure 4).
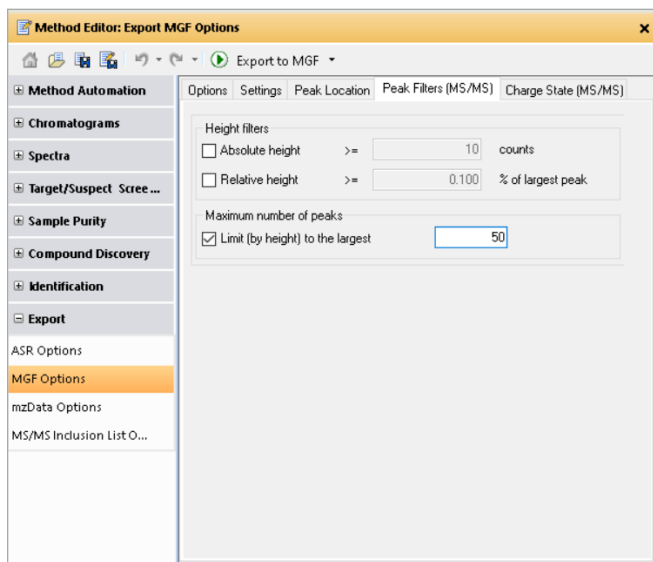


**Figure 4.** Method Editor: Export MGF Options > MGF Options > Peak Filters (MS/MS).

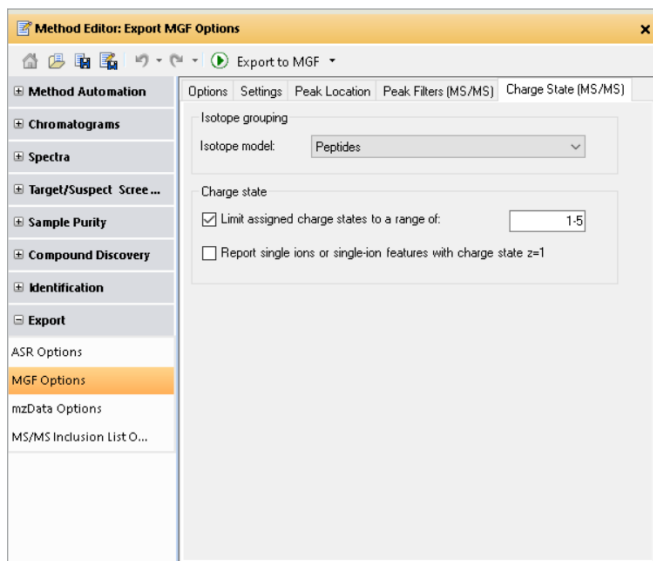7. For **Charge State (MS/MS)**, select **Peptides** as the isotope model and a charge state of 1-5 (Figure 5).



**Figure 5.** Method Editor: Export MGF Options > MGF Options > Charge State (MS/MS).

8. Select the output file location and click **Export to MGF**.

### Data extraction tips and tricks

– Monitor the .MGF file size in Microsoft Windows File Explorer. Similar samples should have similar size files.

– Use an iterative extraction/search approach to optimize parameters.

– Extracting the entire data file will improve the Mascot score, but creates a larger .MGF file.

– Using **Find by AutoMSMS** instead of **Entire Data File** will give a lower Mascot score and a smaller .MGF file.

– Using **Deisotope** improves Mascot scores and reduces the .MGF file size.

– Setting **Maximum Number of Peaks** to 50 will give a higher Mascot score and a smaller .MGF file.

– Setting **Maximum Number of Peaks** to 100 will create an .MGF file approximately twice as large as using 50, but may result in more identifications.

## Using the Mascot MS/MS Ions Search page



**Figure 6.** The Mascot MS/MS Ions Search page.

1. Fill in the information as requested.

2. Select a database. SwissProt is a good start as it is a nonredundant, manually curated, and relatively small database. Because Mascot scoring thresholds are related to database size, larger databases may not increase the number of matches, but will increase search space. An easy way to reduce search space is to select the correct taxonomy within that database. It is good practice to leave the "Contaminate database" in the search.

3. Select the digestion enzyme (this will typically be Trypsin) and allow up to two missed cleavages. Including more than two missed cleavages will increase the search space and slow the process.

4. Quantitation and Crosslinking are advanced features beyond the scope of this technical overview.

5. Fixed modifications would include the alkylating reagent used to block reduced disulfide bonds, usually carbamidomethylation (iodoacetamide) or carboxymethylation (iodoacetic acid). It is recommended to periodically check for incomplete alkylation either by omitting this fixed modification or adding it as a variable modification.

6. Variable modifications as a minimum add oxidation (M) and deamidated (NQ) as the most common. The use of an excessive number of variable modifications greatly increases the search space, slows the process and may increase the FDR.

7. Deselect **Error Tolerant**. This is an advanced feature used to account for artifacts in the data (refer to the Mascot help page for more information).

8. Use Peptide Tol ± 10 ppm MS/MS Tol ± 20 ppm. From a small molecule point of view, these parameters seem excessively wide, but due to the unique way the data are processed, there is little downside and the mass accuracy can be checked in the results. As the peptide tolerance range increases, the search space increases; 20 ppm is the maximum recommended value. The # $^{13}$C is typically set to zero, because the data have been deisotoped as part of data extraction.

9. Select the .MGF data file and upload it to the server center (data will be transferred through the internet). The .MGF data format is the most reliable way to get the extracted data into the Mascot MS/MS Ions Search page. The instrument should be set to ESI-QUAD-TOF. This tells Mascot what type of fragmentation product ions to expect. This cannot be edited in the free online version. If the full product is purchased, you should add the immonium and a-series ions for improved scores.

10. Set Target FDR to 1%, and deselect **Machine Learning**.

11. Click **Start Search**.

The search time online is very fast due to uploading the data to a dedicated server center. The above parameters can be set and saved as a cookie on the web page.[4]

## Mascot MS/MS Ions Search page tips and tricks

– Use the smallest database and taxonomy possible, as this will reduce the search space.

– Try three missed cleavages to test for enzyme digestion completeness.

– Review the reported ppm errors and adjust Peptide Tol and MS/MS Tol, if required.

– Use an iterative approach to optimize parameters.

Results are delivered as a series of web pages, from which results can be copied and pasted (Figures 7 and 8). The protein results are downloadable as a .CSV file from the Report Builder (Figure 9). Graphics, including MS/MS spectra, can be downloaded as an .SVG file. Mascot provides complete online training for its software. Agilent can provide end-to-end proteomic and biopharmaceutical workflows and training.
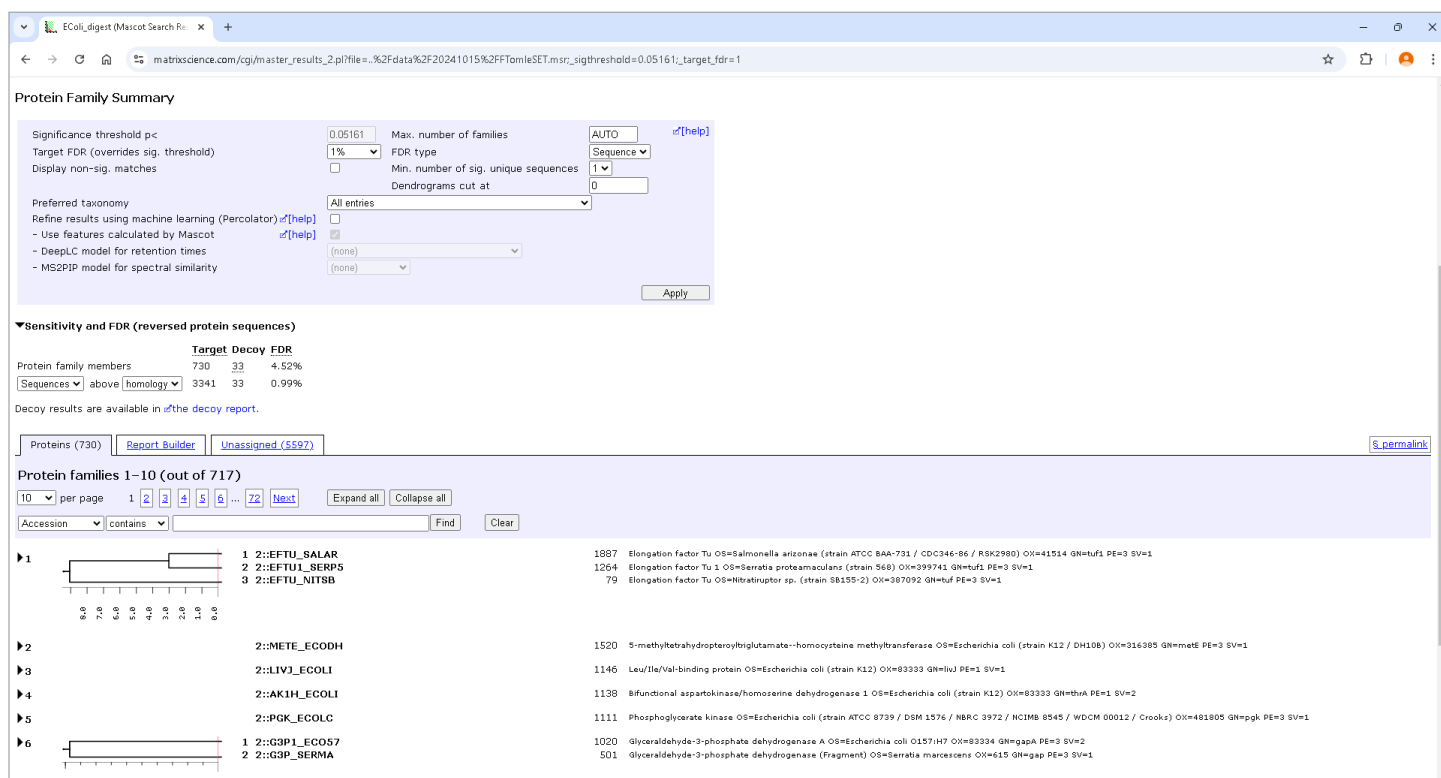


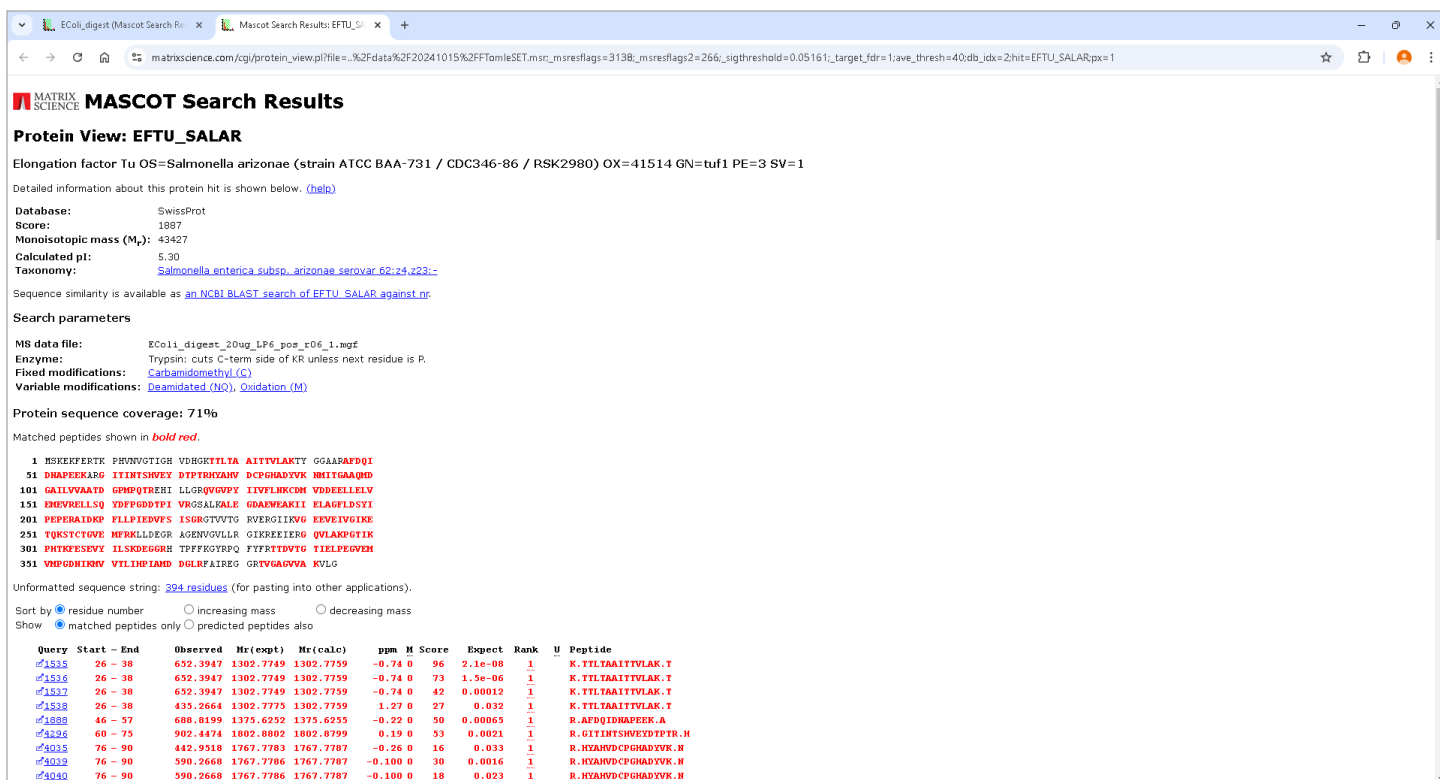**Figure 7.** Mascot search results overview of all proteins.

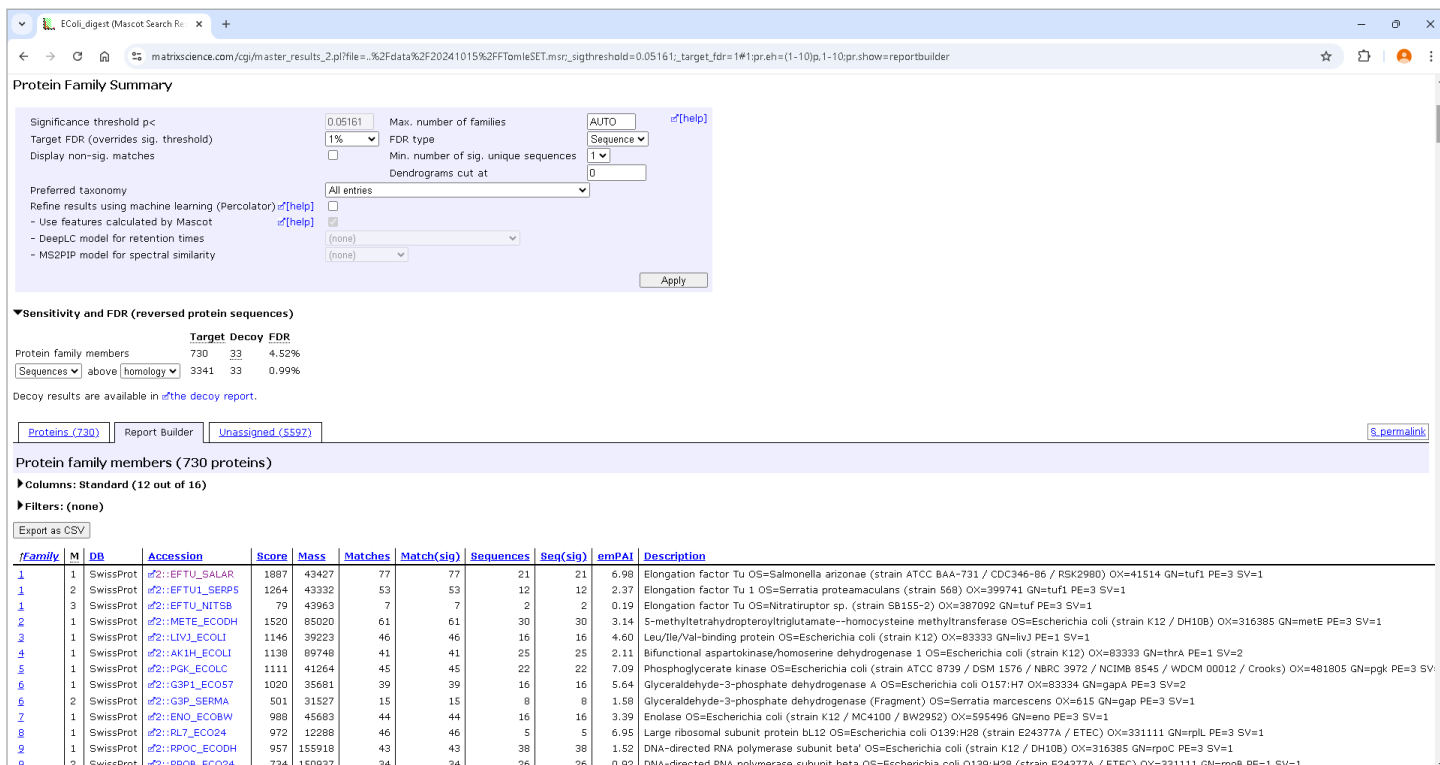**Figure 8.** Mascot search results: single protein details.



**Figure 9.** Report building, export to .CSV file.

## Conclusion

Agilent AutoMSMS proteomic data can be extracted using Agilent MassHunter Qualitative Analysis software, version 10.0 and searched with Mascot online from Matrix Science through an .MGF file. The results can be exported as graphics and as a .CSV file.

## References

1. https://www.matrixscience.com/cgi/search_form.pl?FORMVER=2&SEARCH=MIS

2. https://proteowizard.sourceforge.io/tools/msconvert.html

3. https://www.matrixscience.com/blog/back-to-basics-optimize-your-search-parameters.html#comments

4. https://www.matrixscience.com/search_form_select.html

**www.agilent.com**

Trusted Answers