

# Identifying “Known Unknowns” Using ChemSpider and Automated MS/MS Structure Correlation

James Little, Eastman Chemical Company, Kingsport, TN

Frank Kuhlmann, Xiangdong Li, Jerry Zweigenbaum, Agilent Technologies, Inc., Santa Clara, CA

Antony Williams, Valery Tkachenko, ChemSpider/Royal Society of Chemistry (RSC)

Tuesday Poster Session TP28 Slot: 332

2012 ASMS Conference

60<sup>th</sup> ASMS Conference

Vancouver, Canada

# EASTMAN



**Agilent Technologies**



**RSC** | Advancing the  
Chemical Sciences

## Overview

- “Known Unknowns”<sup>1-3</sup> can be routinely identified using large “spectraless” databases such as ChemSpider with accurate mass MS and MS/MS data
- Known in ChemSpider or other databases, but unknown in the sample to the investigator
- Prototype version of Agilent MassHunter Molecular Structure Correlator (MSC), which is based on a systematic bond-breaking approach,<sup>4</sup> imports associated references in addition to structures when querying ChemSpider by elemental composition
- MSC fragmentation score and assigned substructural information is used in combination with associated references to obtain most probable candidate structures for an unknown

## Introduction

A compound known in the chemical literature is often an unknown to an investigator when being analyzed in a sample. Thus, the origin of the term “known unknown” for this type of non-targeted species.<sup>1-3</sup> Large “spectraless” databases and accurate mass LC-MS data are routinely utilized to identify unknowns prioritized by the number of associated references using these approaches. ChemSpider is a particularly valuable resource since it contains >27 million compounds and is provided as a free resource to the community via the internet.

In previous work, the candidate structures from elemental composition searches were further evaluated by *manually* assigning the observed fragment ions to substructures using a drawing program. The Agilent MassHunter Molecular Structure Correlator (MSC) software *automatically* scores and visually compares the observed MS/MS fragmentation to that of all the candidate structures utilizing a “systematic bond breaking approach.”<sup>4</sup> A unique feature of the MSC program is that it correlates an accurate mass MS/MS spectrum for each candidate structure and assigns a comparative scoring.

In the current work, a prototype MSC program was used which considers the number of associated ChemSpider references together with each structural correlation score to rank possible candidate structures. The resulting structures can then be sorted either by the references or by the structure correlation score. This approach for the identification of unknowns was evaluated with a group of test compounds.

## Methods

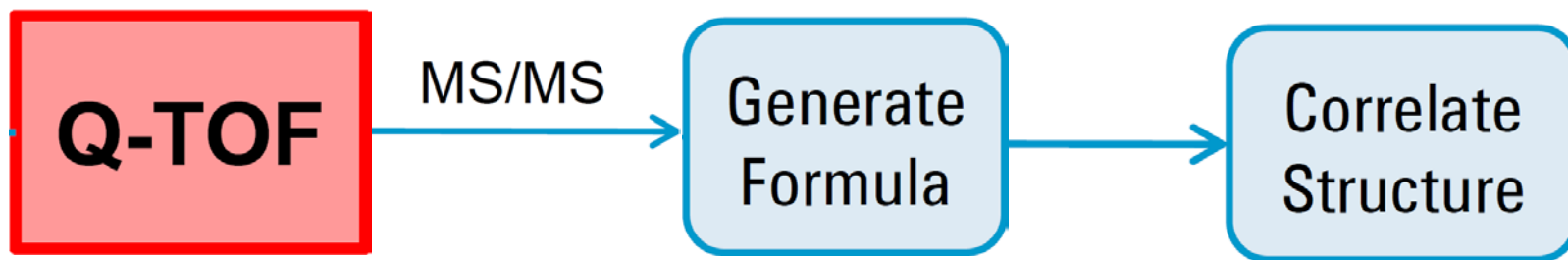
The accurate mass data was obtained for 38 components in Agilent's Toxicology and Environmental Test mixtures. These known compounds were used to evaluate the effectiveness of the approach for the identification of unknowns. The positive ion electrospray accurate mass data was obtained on an Agilent 6540 Q-TOF LC/MS system using Auto MS/MS acquisition.

MassHunter Qualitative Analysis software was used to automatically detect compounds in the test samples using an untargeted approach via the Molecular Feature Extractor (MFE). The latest version of MFE also automatically extracts MS/MS spectra if they exist. For each compound the accurate mass MS and MS/MS information was then exported to a compound exchange file (\*.CEF) and imported into the Agilent Molecular Structure Correlator (MSC) software.

The **current** release version of the MSC software calculates the most probable molecular formula from the accurate mass MS and MS/MS data of an unknown compound and then queries the ChemSpider database via an applications programming interface (API) to automatically retrieve all structures consistent with that elemental composition. The MS/MS spectrum of the unknown compound is then correlated with each candidate structure and an overall correlation score is assigned based on the fraction of explainable fragment ion intensity, their mass accuracy, and the energy needed to create each fragment ion.

The **prototype** MSC program contains additional capabilities to query and process the ChemSpider compound database for the number of associated references. Reference fields and the overall correlation score can then be sorted for evaluation.

## Approach



- Accurate mass data acquired by LC/MS in data-dependent MS/MS mode
- Agilent MassHunter Qualitative Analysis software used to find compounds and generate molecular formulas using monoisotopic mass, isotope abundance, and isotope spacing, as well as look for matching formulas for each fragment ion and its neutral loss from the precursor
- Prototype Agilent MassHunter Molecular Structure Correlator software uses ChemSpider interface to obtain candidate structures for a target molecular formula and associated number of references for candidate structures
- Calculation of correlation scores to explain MS/MS fragmentation pattern for each candidate structure using a “systematic bond-breaking” approach
- Results sorted by number of references and/or correlation score to evaluate structure candidates for identification of components

# Data Processing: MassHunter

Agilent MassHunter Qualitative Analysis 8.05.00 - MSMSLib\_Demo.m

File Edit View Find Identify Chromatograms Spectra Method Sequence Wizards Actions Configuration Tools Help

### Compound List

Automatically Show Columns

| Show/Hide                           | Cpd | Name                   | Score (Li) | RT    | m/z      | Base Pea | Height | Start | End   | Algorithm  | File                  | Z Count |
|-------------------------------------|-----|------------------------|------------|-------|----------|----------|--------|-------|-------|------------|-----------------------|---------|
| <input checked="" type="checkbox"/> | 11  | Methamphetamine        | 99.82      | 2.033 | 150.1277 | 91.0543  | 125194 | 1.941 | 2.126 | Auto MS/MS | For_Tox_mix_1AMSMS2.d | 1       |
| <input checked="" type="checkbox"/> | 12  | Strychnine             | 100        | 2.163 | 335.1756 | 335.1749 | 81256  | 2.163 | 2.163 | Auto MS/MS | For_Tox_mix_1AMSMS2.d | 1       |
| <input checked="" type="checkbox"/> | 13  | Heroin                 | 87.52      | 2.837 | 370.1644 | 370.1629 | 6661   | 2.837 | 2.837 | Auto MS/MS | For_Tox_mix_1AMSMS2.d | 1       |
| <input checked="" type="checkbox"/> | 14  | Cocaine                | 97.25      | 2.947 | 304.1546 | 182.1169 | 78307  | 2.947 | 2.947 | Auto MS/MS | For_Tox_mix_1AMSMS2.d | 1       |
| <input checked="" type="checkbox"/> | 15  | Meperidine (Pethidine) | 96.95      | 3.013 | 248.1651 | 220.1327 | 179735 | 3.013 | 3.013 | Auto MS/MS | For_Tox_mix_1AMSMS2.d | 1       |
| <input checked="" type="checkbox"/> | 17  | Trazodone              | 91.32      | 3.346 | 372.1588 | 176.0809 | 113172 | 3.346 | 3.346 | Auto MS/MS | For_Tox_mix_1AMSMS2.d | 1       |

### Data Navigator

Sort by Data File

- Cpd 1: 0.619
- Cpd 2: 0.795
- Cpd 3: 0.951
- Cpd 4: Codeine
- Cpd 5: Amphetamine
- Cpd 6: Oxycodone
- Cpd 7: 1.830
- Cpd 8: 3,4-Methylenedioxyamphetamine (MDA)
- Cpd 9: Hydrocodone
- Cpd 10: Methylenedioxyamphetamine (MDMA)
- Cpd 11: Methamphetamine
- Cpd 12: Strychnine
- Cpd 13: Heroin
- Cpd 14: Cocaine
- Cpd 15: Meperidine (Pethidine)
- Cpd 17: Trazodone
  - + Scan (3.332 min)
  - + Scan (3.332 min)
  - + Product Ion (3.346 min) (372.1588 -> \*)
  - Cpd 18: 3.506

### Chromatogram Results

Sort by Minutes

x10<sup>5</sup> Cpd 17: Trazodone: +ESI EIC(372.1588) Scan Frag=110.0V For\_Tox\_mix\_1AMSMS2.d

Counts vs. Acquisition Time (min)

### MS Spectrum Results

MS Actuals: + Scan (3.332 min) Sample Information

Chromatogram Spectrum General Reports Find Compounds

Find by Auto MS/MS  
Find by Targeted MS/MS  
Find by Molecular Feature  
Find by Maximum Entropy  
Find by MRM  
Find Compounds by Formula  
Identify Compounds

### MS Spectrum Results

MS Actuals: + Scan (3.332 min) Sample Information

Chromatogram Spectrum General Reports Find Compounds

Find by Auto MS/MS  
Find by Targeted MS/MS  
Find by Molecular Feature  
Find by Maximum Entropy  
Find by MRM  
Find Compounds by Formula  
Identify Compounds

x10<sup>5</sup> Cpd 17: Trazodone: +ESI Scan (3.332 min) Frag=110.0V For\_Tox\_mix\_1AMSMS2.d

Counts vs. Mass-to-Charge (m/z)

x10<sup>4</sup> Cpd 17: Trazodone: +ESI Product Ion (3.346 min) Frag=110.0V CID@20.0 (372.1588[z=1]->\*) For\_Tox\_mix\_1AMSMS2.d

Counts vs. Mass-to-Charge (m/z)

Method Explorer: MSMSLib\_Demo.m Method Editor: Search Accurate Mass Li... Spectrum Preview MS Spectrum Results Spectral Difference Results: Cpd 17: Trazodone

# MSC Prototype Interface

Agilent MassHunter Molecular Structure Correlator B.05.00 -- For\_Tox\_mix\_1AMSMS2\_2--M+H; ce=20

File Settings Help

Compound formula:  Add

M = 371.1515; 4 formula candidates from MFG

| ID | Formula      | Isomers | Taut. Gtps | dM(ppm) | Product | Precur. | Overall |
|----|--------------|---------|------------|---------|---------|---------|---------|
| 1  | C19H22CIN5O  | 296     | 261        | -0.6    | 99      | 98      | 98      |
| 2  | C18H26CINO5  | 97      | 74         | -4.2    | 84      | 91      | 90      |
| 3  | C16H26CIN5OS | 9       | 9          | 8.4     | 99      | 74      | 79      |
| 4  | C22H26CINS   | 10      | 10         | -11.0   | 0       | 64      | 52      |

Fragment formulas for C19H22CIN5O

| #  | m/z      | intensity | nom. intens. | formula    | dM(ppm) |
|----|----------|-----------|--------------|------------|---------|
| 1  | 78.0333  | 1524.64   | 2.85         | C5H4N      |         |
| 2  | 133.0767 | 876.04    | 1.64         | C8H9N2     |         |
| 3  | 133.0767 | 876.04    | 1.64         | C7H14Cl    |         |
| 4  | 148.0500 | 19745.75  | 36.88        | C7H6N3O    |         |
| 5  | 148.0500 | 19745.75  | 36.88        | C4H9CIN4   |         |
| 6  | 148.0500 | 19745.75  | 36.88        | C6H11CINO  |         |
| 7  | 148.0857 | 3728.19   | 6.96         | C8H10N3    |         |
| 8  | 148.0857 | 3728.19   | 6.96         | C10H12O    |         |
| 9  | 148.0857 | 3728.19   | 6.96         | C7H15CIN   |         |
| 10 | 176.0809 | 53535.29  | 100.00       | C9H10N3O   |         |
| 11 | 176.0809 | 53535.29  | 100.00       | C6H13CIN4  |         |
| 12 | 176.0809 | 53535.29  | 100.00       | C8H15CINO  |         |
| 13 | 177.0823 | 640.31    | 1.20         | C7H14CIN2O |         |
| 14 | 177.0823 | 640.31    | 1.20         | C5H12CIN5  |         |
| 15 | 177.0823 | 640.31    | 1.20         | C8H9N4O    |         |
| 16 | 209.0815 | 316.04    | 0.59         | C12H9N4    |         |
| 17 | 209.0815 | 316.04    | 0.59         | C14H11NO   |         |
| 18 | 209.0815 | 316.04    | 0.59         | C11H14CIN2 |         |
| 19 | 235.0928 | 276.16    | 0.52         | C14H16CIO  |         |
| 20 | 235.0928 | 276.16    | 0.52         | C14H11N4   |         |
| 21 | 235.0928 | 276.16    | 0.52         | C12H14CIN3 |         |
| 22 | 235.0928 | 276.16    | 0.52         | C15H11N2O  |         |

Compound formula: C19H22CIN5O

Structure Search Parameters: Compatibles/Total: 187/244

ChemSpider (Web) Go

Add Structure

Sort by: Score

Score

# Data sources

# References

# PubMed

# RSC

Standard InChIKey: 1  
PHLBKPHSAVXXEF-UHFFFAOYSA-N  
Score: 92.88  
More Info ...  
MSC Save Delete  
ChemSpider: 5332

ChemSpider Info

# of References=1007

# of Data sources=41

# of PubMed=1325

# of RSC=9

Standard InChIKey: 2  
WZFUFXDNTAVTTF-UHFFFAOYSA-N  
Score: 92.88  
More Info ...  
MSC Save Delete  
ChemSpider: 7977308

Standard InChIKey: 3  
MWVNIHSTAKQEOP-UHFFFAOYSA-N  
Score: 91.91  
More Info ...

Fragments of structure #1 -- elucidated: Display Filters

| # | Mass     | Intensity | Weight(%) | No. of candid. | Best score |
|---|----------|-----------|-----------|----------------|------------|
| 1 | 176.0809 | 53535.29  | 75.9      | 2              | 95.7       |
| 2 | 148.0500 | 19745.75  | 9.9       | 7              | 96.7       |
| 3 | 148.0857 | 3728.19   | 1.9       | 6              | 90.7       |
| 4 | 78.0333  | 1524.64   | 0.0       | 3              | 88.4       |
| 5 | 237.1147 | 912.79    | 7.7       | 1              | 97.5       |
| 6 | 133.0767 | 876.04    | 0.2       | 3              | 79.8       |
| 7 | 177.0823 | 640.31    | 0.9       | 6              | 28.3       |
| 8 | 209.0815 | 316.04    | 1.3       | 3              | 80.6       |
| 9 | 235.0928 | 276.16    | 2.2       | 0              | 0.0        |

Penalty=2.0 dM=5.3ppm Score=95.7 Of 1 Penalty=8.5 dM  
C9H11N3O-H C9H15N3O-5H

Chemical structures are displayed for each entry, with red numbers 1-8 indicating specific features or annotations.

## Windows in MSC Prototype

- 1.** Molecular formulas with MS score (mass position, isotope abundance, isotope spacing), MS/MS score and combined score
- 2.** Molecular formulas for fragment ions with  $m/z$  error, part of MS/MS score
- 3.** Number of structures retrieved for formula/number of structures qualified in structure correlation
- 4.** Fields used to sort the results including reference categories from ChemSpider and overall structure correlation score calculated by MSC program
- 5.** Structure selected for showing substructural information in sections **7** and **8**
- 6.** Number of references retrieved from ChemSpider for each category
- 7.** Tabulated information for substructures including best individual fragment score
- 8.** Substructures for selected structure with score, mass error, and penalty



# Candidate Structure Filters in MSC

Search Parameters

Min Compatibility Score

Max no. of compounds

Representative tautomer only

No. of total rings  
 unconstrained  
Min   
Max

No. of Aromatic rings  
 unconstrained  
Min   
Max

Filter by ChemSpider Meta Data

Min # of references

Min # of data sources

Min # of PubMed

Min # of RSC

Functional-Group Filters

Display Filters

Fragment peak filter

None  
 By abs. height  
 By rel. height  
 By weight

Min

Show isotopes  
 Hide  
 Show

Substructure candidate filter

High-quality candidate definition

By abs. score  %  
 By rel. score

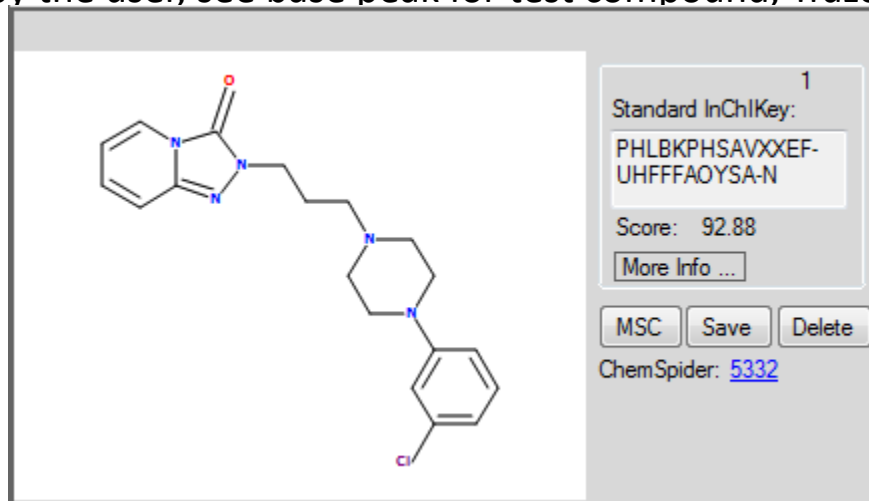
Low-quality candidate truncation

Max total candidate count

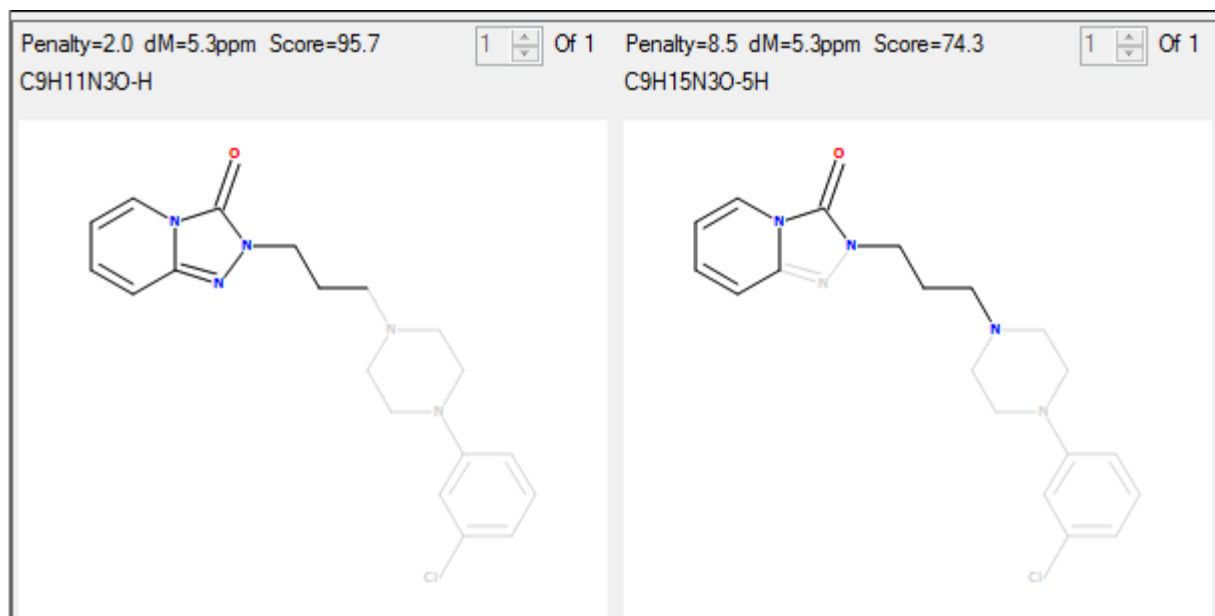
# Easy Review of Substructures for Fragment Ions

The proposed substructures and associated individual fragment ion score are easily reviewed by the user, see base peak for test compound, Trazodone, in study

5



8



# Partial Results of Searches

| Compound                      | No. Structure Candidates | Structure Score | Structure Rank | No. References | Reference Rank |
|-------------------------------|--------------------------|-----------------|----------------|----------------|----------------|
| ecgonine methyl ester         | 743                      | 87.21           | 7              | 156            | 1              |
| codeine                       | 2769                     | 73.6            | 104            | 2246           | 1              |
| amphetamine                   | 277                      | 94.5            | 3              | 10829          | 1              |
| oxycodone                     | 2371                     | 87.07           | 94             | 611            | 1              |
| 3,4-methylenedioxyamphetamine | 1494                     | 92.26           | 32             | 272            | 2              |
| hydrocodone                   | 2769                     | 66.33           | 212            | 201            | 2              |
| methamphetamine               | 446                      | 95.91           | 26             | 3559           | 1              |
| strychnine                    | 2256                     | 14.58           | 195            | 2468           | 1              |
| cocaine                       | 1417                     | 95.44           | 9              | 15414          | 1              |
| meperidine                    | 1953                     | 97.46           | 27             | 2959           | 1              |
| trazodone                     | 244                      | 93.17           | 1              | 956            | 1              |
| phencyclidine                 | 274                      | 94.01           | 12             | 2524           | 1              |
| verapamil                     | 210                      | 93.96           | 4              | 13379          | 1              |
| levomethadone                 | 758                      | 89.93           | 13             | 5713           | 1              |
| nitrazepam                    | 568                      | 86.07           | 49             | 765            | 1              |
| alprazolam xanax              | 104                      | 81.71           | 16             | 1229           | 1              |
| clonazepam                    | 210                      | 87.62           | 18             | 1975           | 1              |
| proadifen                     | 594                      | 92.39           | 8              | 272            | 1              |
| temazepam                     | 786                      | 94.14           | 27             | 1164           | 1              |
| diazepam                      | 494                      | 89.03           | 10             | 12751          | 1              |
| diethyl phthalate             | 1264                     | 95.35           | 16             | 179            | 1              |
| aminocarb                     | 1163                     | 97.37           | 7              | 97             | 3              |
| thiabendazole                 | 75                       | 95.52           | 18             | 1422           | 1              |
| imazapyr                      | 1040                     | 94.61           | 5              | 88             | 1              |
| dimethoate                    | 2                        | 94.89           | 1              | 660            | 1              |
| imazapic                      | 1322                     | 88.82           | 16             | 30             | 1              |
| metoxuron                     | 205                      | 88.76           | 7              | 67             | 1              |
| enilconazole                  | 104                      | 81.42           | 14             | 171            | 1              |
| atrazine                      | 22                       | 97.24           | 2              | 1776           | 1              |
| metosulam                     | 11                       | 43.93           | 7              | 25             | 1              |
| metazachlor                   | 440                      | 96.42           | 1              | 46             | 1              |
| pyraclostrobin                | 840                      | 85.51           | 1              | 19             | 1              |
| monocrotaline                 | 303                      | 39.42           | 184            | 556            | 1              |
| simazine                      | 15                       | 82.16           | 7              | 522            | 1              |
| molinate                      | 99                       | 92.38           | 6              | 153            | 1              |
| malathion                     | 8                        | 91.68           | 4              | 1794           | 1              |
| heroin                        | 2680                     | 85.94           | 4              | 147            | 1              |

# Tabulated Results

|                                  |                                 |                                   |   |
|----------------------------------|---------------------------------|-----------------------------------|---|
| <b>No. Candidates Average</b>    | <b>No. Candidates Median</b>    | <b>Range of Candidates</b>        |   |
| 849                              | 531                             | 22-769                            |   |
| <b>Average Structure Score</b>   | <b>Median Structure Score</b>   | <b>Range Structure Score</b>      |   |
| 85                               | 91                              | 15-97                             |   |
| <b>Average Structure Rank</b>    | <b>Median Structure Rank</b>    | <b>Range of Structure Rank</b>    |   |
| 32                               | 11                              | 1-212                             |   |
| <b>Average No. of References</b> | <b>Median No. of References</b> | <b>Range of References</b>        |   |
| 2319                             | 713                             | 13-15000                          |   |
| <b>Average Reference Rank</b>    | <b>Median Reference Rank</b>    | <b>Range of Rank</b>              | <b>Entries with No Rank (0 Entries)</b> |
| 1                                | 1                               | 1-3                               | 0                                       |
| <b>Average PubMed References</b> | <b>Median PubMed References</b> | <b>Range of PubMed References</b> |   |
| 2328                             | 24                              | 0-26,000                          |   |
| <b>Average PubMed Rank</b>       | <b>Median PubMed Rank</b>       | <b>Range of PubMed Rank</b>       | <b>Entries with No Rank (0 Entries)</b> |
| 2                                | 1                               | 1-7                               | 9                                       |
| <b>Average RSC References</b>    | <b>Median RSC References</b>    | <b>Range of RSC References</b>    |   |
| 245                              | 4                               | 0-1737                            |   |
| <b>Average RSC Rank</b>          | <b>Median RSC Rank</b>          | <b>Range of RSC Rank</b>          | <b>Entries with No Rank (0 Entries)</b> |
| 2                                | 1                               | 1-7                               | 11                                      |
| <b>Average Data Sources</b>      | <b>Median Data Sources</b>      | <b>Range of Data Sources</b>      |   |
| 27                               | 24                              | 12-63                             |   |
| <b>Average RSC Rank</b>          | <b>Median RSC Rank</b>          | <b>Range of RSC Rank</b>          | <b>Entries with No Rank (0 Entries)</b> |
| 3                                | 1                               | 1-7                               | 0                                       |

## Conclusions

- Initial results with the limited data set indicate that the MassHunter MSC program is a very useful tool for the identification of unknown compounds in mixtures especially with the addition of the associated reference fields
- Hundreds of candidate structures can be returned from a large compound database like ChemSpider for one elemental composition; thus, there can be many hits with a similar correlation score (and similar chemical structure) which makes it difficult to pinpoint the correct structure for the unknown compound
- Sorting of the candidate structures by “number of references” is most useful for bringing the proper identification to the top of the hit list
- Sorting by PubMed or RSC references and data sources also proves to increase the discrimination between the hits
- The combination of the MSC correlation score in combination with the number of references results in getting the correct hits to the top of the results list
- Efficient substructure viewing in MSC greatly helps in validating the *most likely* candidate after first sorting by “number of references” field and then the correlation score

## Possible Future Work

- Implement the ability to import and sort by references into the next version of MSC supplied with MassHunter and provide the meta information about the number of references in the ChemSpider application programming interface
- Provide a weighted combined scoring between the correlation score and the different no. of references criteria to increase rank of correct structure for known test compounds
- Refine MSC interface for efficient sorting/displaying of ChemSpider reference fields

## Acknowledgments

Many thanks to Curt Cleven, Jean Coffman, Mike Ramsey, Stacy Brown, Alexey Pshenichnov, Mike Scott, Bill Tindall, and Kent Morrill for their contributions to the various approaches for identifying “known unknowns.”

## References

1. Identification of “Known Unknowns” Utilizing Accurate Mass Data and Chemical Abstracts Service, J. Little, C. Cleven, and S. Brown, JASMS, Vol 22, No. 2 (2011), 348-359.
2. Identification of “Known Unknowns” Utilizing Accurate Mass Data and ChemSpider, J. Little, A. Williams, A. Pshenichnov, V. Tkachenko, JASMS, Vol 23, No. (2012), 179-185.
3. [Littlemsandsailing.wordpress.com](http://Littlemsandsailing.wordpress.com)
4. Automated Assignment of High-Resolution Collisionally Activated Dissociation Mass Spectra Using a Systematic Bond Disconnection Approach, A. W. Hill, R. J. Mortishire-Smith, Rapid Communications in Mass Spectrometry, Vol 19, Issue 21 (2005) 3111-3118.