

**This item is the archived peer-reviewed author-version of:**

spectrum\_utils : a python package for mass spectrometry data processing and visualization

**Reference:**

Bittremieux Wout.- spectrum\_utils : a python package for mass spectrometry data processing and visualization  
Analytical chemistry - ISSN 0003-2700 - 92:1(2020), p. 659-661  
Full text (Publisher's DOI): <https://doi.org/10.1021/ACS.ANALCHEM.9B04884>  
To cite this reference: <https://hdl.handle.net/10067/1656830151162165141>

## Technical Note

spectrum\_utils: A Python package for mass  
spectrometry data processing and visualization

Wout Bittremieux

*Anal. Chem.*, **Just Accepted Manuscript** • DOI: 10.1021/acs.analchem.9b04884 • Publication Date (Web): 06 Dec 2019

Downloaded from pubs.acs.org on December 6, 2019

**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 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 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.

# spectrum\_utils: A Python package for mass spectrometry data processing and visualization

Wout Bittremieux<sup>1,2,3,\*</sup>

<sup>1</sup>Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, La Jolla, CA 92093, USA; <sup>2</sup>Department of Mathematics and Computer Science, University of Antwerp, 2020 Antwerp, Belgium; <sup>3</sup>Biomedical Informatics Network Antwerpen (biomina), 2020 Antwerp, Belgium.

Corresponding author: [wbittremieux@health.ucsd.edu](mailto:wbittremieux@health.ucsd.edu)

## Abstract

Given the wide diversity in applications of biological mass spectrometry, custom data analyses are often needed to fully interpret the results of an experiment. Such bioinformatics scripts necessarily include similar basic functionality to read mass spectral data from standard file formats, process it, and visualize it. Rather than having to reimplement this functionality, to facilitate this task, `spectrum_utils` is a Python package for mass spectrometry data processing and visualization. Its high-level functionality enables developers to quickly prototype ideas for computational mass spectrometry projects in only a few lines of code. Notably, the data processing functionality is highly optimized for computational efficiency to be able to deal with the large volumes of data that are generated during mass spectrometry experiments. The visualization functionality makes it possible to easily produce publication-quality figures as well as interactive spectrum plots for inclusion on web pages.

`spectrum_utils` is available for Python 3.6+, includes extensive online documentation and examples, and can be easily installed using conda. It is freely available as open source under the Apache 2.0 license at [https://github.com/bittremieux/spectrum\\_utils](https://github.com/bittremieux/spectrum_utils).

## Introduction

Mass spectrometry (MS) is a powerful, high-throughput analytical technique that can be used to identify and quantify molecules in complex biological samples. Because during a typical MS experiment tens of thousands of mass spectra are generated, suitable bioinformatics tools are needed to analyze such large data volumes. MS data processing has traditionally been done using monolithic software tools that aim to provide fully end-to-

1  
2 end solutions from the raw data to the final identification or quantification results. However, because there  
3 exist a large variety of experimental set-ups and configurations, such tools necessarily cannot cover all possible  
4 use cases.  
5  
6  
7

8  
9 Instead, customized data analysis workflows are often needed to fully interpret the results of an MS experi-  
10 ment. In recent years several software packages for the general-purpose analysis of MS data in popular script-  
11 ing languages have been developed. Notable examples include MSnbase<sup>1</sup> for data processing, visualization,  
12 and quantification in R; pymzML<sup>2,3</sup> to efficiently read and process spectra in the mzML format<sup>4</sup> using Python;  
13 Pyteomics<sup>5,6</sup> for a variety of proteomics data processing tasks in Python; and pyOpenMS<sup>7</sup> to expose the rich  
14 functionality of OpenMS<sup>8,9</sup> from C++ to Python.  
15  
16  
17  
18  
19

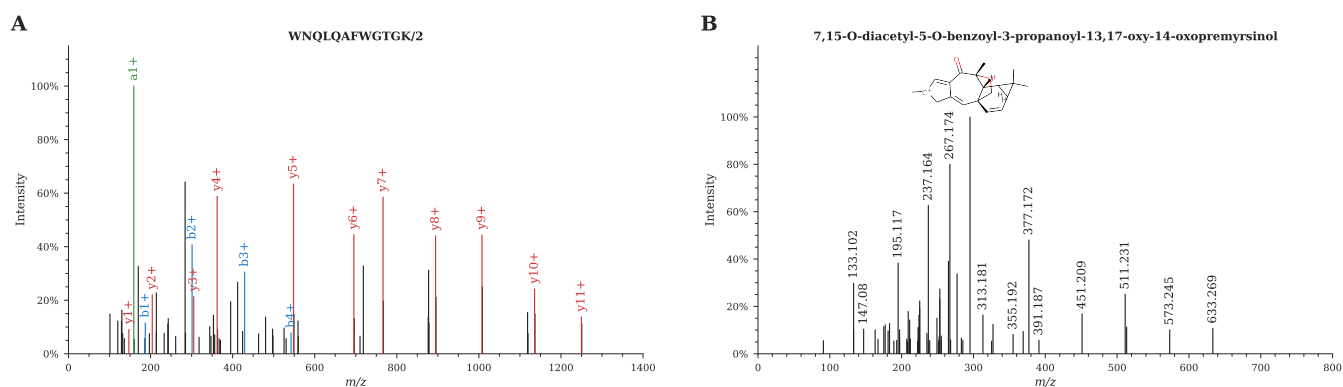
20 Here we present the `spectrum_utils` package for MS data processing and visualization in Python.  
21 `spectrum_utils` allows the user to easily manipulate mass spectral data and quickly prototype ideas for com-  
22 putational MS projects. A key feature of `spectrum_utils` is its focus on computational efficiency to process  
23 large amounts of spectral data. `spectrum_utils` is freely available as open source under the Apache 2.0 license  
24 at [https://github.com/bittremieux/spectrum\\_utils](https://github.com/bittremieux/spectrum_utils).  
25  
26  
27  
28  
29  
30  
31  
32

## 33 **Methods & results**

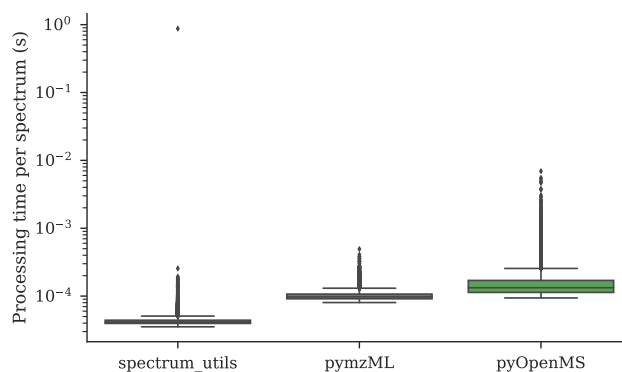
34  
35  
36  
37

38 The functionality provided by `spectrum_utils` is built around the concept of tandem mass spectrometry  
39 (MS/MS) spectra as basic elements. MS/MS spectra can be processed in various ways. Uninformative peaks,  
40 such as the precursor ion and its isotopic peaks, and low-intensity noise peaks, can be removed. Further peak  
41 filtering can be performed by only retaining the top most intense peaks. Next, peak intensities can be scaled  
42 to de-emphasize overly dominant peaks. Possible transformations are root scaling, log scaling, or rank-based  
43 scaling. Finally, peaks can be annotated with their peptide fragments, potentially including post-translational  
44 modifications (PTMs) at specified amino acid positions, molecules encoded as SMILES strings, or custom an-  
45 notations (figure 1).  
46  
47  
48  
49  
50  
51  
52

53 An important emphasis is placed on the computational efficiency of the spectrum processing steps. Operations  
54 on the peaks of MS/MS spectra are implemented using NumPy,<sup>12</sup> a popular Python library for efficient numeri-  
55 cal computation. Additionally, Numba,<sup>13</sup> a Python just-in-time (JIT) compiler, is used to further speed up many  
56 processing steps by compiling Python and NumPy code into efficient machine instructions. Using these tech-  
57 niques for optimized numerical computation, `spectrum_utils` is able to very efficiently process large amounts  
58 of MS data (figure 2).  
59  
60



**Figure 1:** Visualization of (A) a proteomics MS/MS spectrum (mzspec:PXD004732:01650b\_BC2-TUM\_first\_pool\_53\_01\_01-3xHCD-1h-R2:scan:41840)<sup>10</sup> and (B) a metabolomics MS/MS spectrum (mzspec:GNPSLIBRARY:CCMSLIB00000840351).<sup>11</sup> Fragment peaks can be easily annotated based on a peptide sequence, a SMILES string corresponding to a (sub)structure, or using custom annotations.



**Figure 2:** Spectrum processing runtime comparison. The runtime for processing the spectra generated for the iPRG2012 challenge<sup>14</sup> is reported for similar processing steps for spectrum\_utils (version 0.3.0), pymzML (version 2.4.4),<sup>2,3</sup> and pyOpenMS (version 2.4.0).<sup>7</sup> The corresponding Python script is available in the supplementary information and at the online spectrum\_utils documentation. Spectrum processing consisted of fixing the  $m/z$  range, removing precursor ion peaks and low-intensity noise peaks, and scaling the peak intensities by their square root. Note that the significant outlier for spectrum\_utils is caused by Numba's JIT compilation<sup>13</sup> of the first method call, allowing subsequent calls to be made very efficiently.

MS/MS spectra can be easily visualized using spectrum\_utils's plotting functionality to help in producing publication-quality figures (figure 1). Additionally, a mirror plot can be used to clearly show the similarity between two different spectra, for example, to visualize identification results produced by spectral library searching. The default plotting functionality uses matplotlib<sup>15</sup> to generate high-quality spectrum figures. Alternatively, interactive plotting functionality is available using Altair,<sup>16</sup> which is based on the Vega and Vega-Lite grammar of interactive graphics.<sup>17</sup> Interactive plotting is a drop-in replacement for the standard plotting functionality, allowing the user to trivially produce interactive spectrum plots for data exploration or visualization on web pages.

## Use cases

We will briefly highlight two use cases to demonstrate the functionality of `spectrum_utils`. First, `spectrum_utils` is used by the ANN-SoLo open modification spectral library search engine<sup>18,19</sup> to preprocess MS/MS spectra prior to spectral library searching. `spectrum_utils` provides the full functionality needed to preprocess the spectral data, including rounding the peak  $m/z$  values to a common mass resolution, removing the precursor peak and low-intensity noise peaks, and scaling the peak intensities. ANN-SoLo was developed to efficiently perform open modification searches, which typically suffer from a very large search space. In part due to the efficient spectrum processing functionality provided by `spectrum_utils`, ANN-SoLo is able to identify hundreds to thousands of spectra per minute, enabling researchers to perform untargeted PTM profiling via open searches at an unprecedented scale.

Second, the Global Natural Products Social Molecular Networking (GNPS) public data repository and analysis platform uses `spectrum_utils` to produce high-quality spectrum vector graphics. All MS/MS spectra processed by GNPS can be visualized using their universal spectrum identifier,<sup>20</sup> as well as the individual spectra in reference spectral libraries compiled using the MSMS-Chooser workflow.<sup>21</sup> Additionally, spectral library identification results obtained using GNPS are visualized using mirror plots.

## Code availability

`spectrum_utils` is available for Python 3.6+ and can be easily installed via conda using the Bioconda channel.<sup>22</sup> `spectrum_utils` depends on Numpy<sup>12</sup> and Numba<sup>13</sup> for efficient numerical computation, Pyteomics<sup>5,6</sup> for peptide fragment ion mass calculations, RDKit<sup>23</sup> for SMILES string handling, matplotlib<sup>15</sup> for static plotting, and Altair<sup>16</sup> and Pandas<sup>24</sup> for interactive plotting.

All code and detailed documentation on how to use `spectrum_utils` is freely available as open source under the Apache 2.0 license at [https://github.com/bittremieux/spectrum\\_utils](https://github.com/bittremieux/spectrum_utils).

## Conclusion

Here we have presented the `spectrum_utils` package for MS data processing and visualization in Python. Its clearly defined functionality allows `spectrum_utils` to fill an important gap in the Python MS processing ecosystem. For example, `spectrum_utils` does not provide functionality to read spectral data files because

1  
2 there already exist several excellent tools to perform this task. Instead, `spectrum_utils` takes the MS data pro-  
3 vided by such tools as input for subsequent processing and visualization. `spectrum_utils` has a well-defined,  
4 Pythonic application programming interface, allowing developers to easily harness its powerful functionality  
5 in a small number of lines of code. Additionally, it is trivial to switch between the static plotting functionality  
6 to produce high-quality spectrum figures to include in scientific manuscripts and the interactive plotting func-  
7 tionality. This interactive plotting functionality can be very powerful during an explorative analysis of MS data  
8 or to include dynamic spectrum plots in web resources.  
9  
10  
11  
12  
13  
14

15 Besides the two use cases highlighted above, `spectrum_utils` has already been used in several other compu-  
16 tational MS projects, showcasing its versatile functionality. Among others, it has been used to perform custom  
17 analyses of MS data corresponding to an unknown protein sample,<sup>25</sup> to prepare MS/MS spectra for processing  
18 using deep neural networks,<sup>26</sup> and to visualize spectra for scientific publications and presentations. This illus-  
19 trates how `spectrum_utils` facilitates several MS data processing and visualization tasks and allows developers  
20 to quickly prototype ideas and implement diverse computational MS projects.  
21  
22  
23  
24  
25  
26  
27  
28

## 29 Supporting information

30  
31  
32  
33 Benchmark script to compare the spectrum processing runtime between `spectrum_utils`, `pymzML` and `py-`  
34 `OpenMS`.  
35  
36  
37  
38  
39

## 40 Funding

41  
42  
43  
44 W.B. is a postdoctoral researcher of the Research Foundation – Flanders (FWO).  
45  
46  
47  
48  
49

## 50 References

- 51  
52  
53 (1) Gatto, L., Lilley, K. S. MSnbase-an R/Bioconductor Package for Isobaric Tagged Mass Spectrometry  
54 Data Visualization, Processing and Quantitation. *Bioinformatics* **2012**, *28*, 288–289, DOI: [10.1093 /  
55 bioinformatics/btr645](https://doi.org/10.1093/bioinformatics/btr645).  
56  
57  
58  
59  
60

- (2) Bald, T., Barth, J., Niehues, A., Specht, M., Hippler, M., Fufezan, C. pymzML—Python Module for High-Throughput Bioinformatics on Mass Spectrometry Data. *Bioinformatics* **2012**, *28*, 1052–1053, DOI: [10.1093/bioinformatics/bts066](https://doi.org/10.1093/bioinformatics/bts066).
- (3) Kösters, M., Leufken, J., Schulze, S., Sugimoto, K., Klein, J. A., Zahedi, R. P., Hippler, M., Leidel, S. A., Fufezan, C. pymzML v2.0: Introducing a Highly Compressed and Seekable Gzip Format. *Bioinformatics* **2018**, *34*, ed. by Wren, J., 2513–2514, DOI: [10.1093/bioinformatics/bty046](https://doi.org/10.1093/bioinformatics/bty046).
- (4) Martens, L., Chambers, M., Sturm, M., Kessner, D., Levander, F., Shofstahl, J., Tang, W. H., Römpf, A., Neumann, S., Pizarro, A. D., et al. mzML—a Community Standard for Mass Spectrometry Data. *Molecular & Cellular Proteomics* **2011**, *10*, R110.000133–R110.000133, DOI: [10.1074/mcp.R110.000133](https://doi.org/10.1074/mcp.R110.000133).
- (5) Goloborodko, A. A., Levitsky, L. I., Ivanov, M. V., Gorshkov, M. V. Pyteomics—a Python Framework for Exploratory Data Analysis and Rapid Software Prototyping in Proteomics. *Journal of The American Society for Mass Spectrometry* **2013**, *24*, 301–304, DOI: [10.1007/s13361-012-0516-6](https://doi.org/10.1007/s13361-012-0516-6).
- (6) Levitsky, L. I., Klein, J. A., Ivanov, M. V., Gorshkov, M. V. Pyteomics 4.0: Five Years of Development of a Python Proteomics Framework. *Journal of Proteome Research* **2019**, *18*, 709–714, DOI: [10.1021/acs.jproteome.8b00717](https://doi.org/10.1021/acs.jproteome.8b00717).
- (7) Röst, H. L., Schmitt, U., Aebersold, R., Malmström, L. pyOpenMS: A Python-Based Interface to the OpenMS Mass-Spectrometry Algorithm Library. *PROTEOMICS* **2014**, *14*, 74–77, DOI: [10.1002/pmic.201300246](https://doi.org/10.1002/pmic.201300246).
- (8) Sturm, M., Bertsch, A., Gröpl, C., Hildebrandt, A., Hussong, R., Lange, E., Pfeifer, N., Schulz-Trieglaff, O., Zerck, A., Reinert, K., et al. OpenMS – An Open-Source Software Framework for Mass Spectrometry. *BMC Bioinformatics* **2008**, *9*, 163, DOI: [10.1186/1471-2105-9-163](https://doi.org/10.1186/1471-2105-9-163).
- (9) Röst, H. L., Sachsenberg, T., Aiche, S., Bielow, C., Weisser, H., Aicheler, F., Andreotti, S., Ehrlich, H.-C., Gutenbrunner, P., Kenar, E., et al. OpenMS: A Flexible Open-Source Software Platform for Mass Spectrometry Data Analysis. *Nature Methods* **2016**, *13*, 741–748, DOI: [10.1038/nmeth.3959](https://doi.org/10.1038/nmeth.3959).
- (10) Zolg, D. P., Wilhelm, M., Schnatbaum, K., Zerweck, J., Knaute, T., Delanghe, B., Bailey, D. J., Gessulat, S., Ehrlich, H.-C., Weininger, M., et al. Building ProteomeTools Based on a Complete Synthetic Human Proteome. *Nature Methods* **2017**, DOI: [10.1038/nmeth.4153](https://doi.org/10.1038/nmeth.4153).
- (11) Nothias-Esposito, M., Nothias, L. F., Da Silva, R. R., Retailleau, P., Zhang, Z., Leyssen, P., Roussi, F., Touboul, D., Paolini, J., Dorrestein, P. C., et al. Investigation of Premyrsinane and Myrsinane Esters in *Euphorbia Cupanii* and *Euphorbia Pithyusa* with MS2LDA and Combinatorial Molecular Network Annotation Propagation. *Journal of Natural Products* **2019**, *82*, 1459–1470, DOI: [10.1021/acs.jnatprod.8b00916](https://doi.org/10.1021/acs.jnatprod.8b00916).
- (12) Van der Walt, S., Colbert, S. C., Varoquaux, G. The NumPy Array: A Structure for Efficient Numerical Computation. *Computing in Science & Engineering* **2011**, *13*, 22–30, DOI: [10.1109/MCSE.2011.37](https://doi.org/10.1109/MCSE.2011.37).
- (13) Lam, S. K., Pitrou, A., Seibert, S. In *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC - LLVM '15*, ACM Press: Austin, TX, USA, 2015, pp 1–6, DOI: [10.1145/2833157.2833162](https://doi.org/10.1145/2833157.2833162).



- 1  
2 (14) Chalkley, R. J., Bandeira, N., Chambers, M. C., Clauser, K. R., Cottrell, J. S., Deutsch, E. W., Kapp,  
3 E. A., Lam, H. H. N., McDonald, W. H., Neubert, T. A., et al. Proteome Informatics Research Group  
4 (iPRG)\_2012: A Study on Detecting Modified Peptides in a Complex Mixture. *Molecular & Cellular Pro-*  
5 *teomics* **2014**, *13*, 360–371, DOI: [10.1074/mcp.M113.032813](https://doi.org/10.1074/mcp.M113.032813).  
6  
7  
8 (15) Hunter, J. D. Matplotlib: A 2D Graphics Environment. *Computing in Science & Engineering* **2007**, *9*, 90–95,  
9 DOI: [10.1109/MCSE.2007.55](https://doi.org/10.1109/MCSE.2007.55).  
10  
11 (16) VanderPlas, J., Granger, B., Heer, J., Moritz, D., Wongsuphasawat, K., Satyanarayan, A., Lees, E., Timo-  
12 feev, I., Welsh, B., Sievert, S. Altair: Interactive Statistical Visualizations for Python. *Journal of Open Source*  
13 *Software* **2018**, *3*, 1057, DOI: [10.21105/joss.01057](https://doi.org/10.21105/joss.01057).  
14  
15 (17) Satyanarayan, A., Moritz, D., Wongsuphasawat, K., Heer, J. Vega-Lite: A Grammar of Interactive Graph-  
16 ics. *IEEE Transactions on Visualization and Computer Graphics* **2016**, *23*, 341–350, DOI: [10.1109/TVCG.2016.](https://doi.org/10.1109/TVCG.2016.2599030)  
17 [2599030](https://doi.org/10.1109/TVCG.2016.2599030).  
18  
19 (18) Bittremieux, W., Meysman, P., Noble, W. S., Laukens, K. Fast Open Modification Spectral Library Search-  
20 ing through Approximate Nearest Neighbor Indexing. *Journal of Proteome Research* **2018**, *17*, 3463–3474,  
21 DOI: [10.1021/acs.jproteome.8b00359](https://doi.org/10.1021/acs.jproteome.8b00359).  
22  
23 (19) Bittremieux, W., Laukens, K., Noble, W. S. Extremely Fast and Accurate Open Modification Spectral  
24 Library Searching of High-Resolution Mass Spectra Using Feature Hashing and Graphics Processing  
25 Units. *Journal of Proteome Research* **2019**, *18*, 3792–3799, DOI: [10.1021/acs.jproteome.9b00291](https://doi.org/10.1021/acs.jproteome.9b00291).  
26  
27 (20) Deutsch, E. W., Orchard, S., Binz, P.-A., Bittremieux, W., Eisenacher, M., Hermjakob, H., Kawano, S., Lam,  
28 H., Mayer, G., Menschaert, G., et al. Proteomics Standards Initiative: Fifteen Years of Progress and Future  
29 Work. *Journal of Proteome Research* **2017**, *16*, 4288–4298, DOI: [10.1021/acs.jproteome.7b00370](https://doi.org/10.1021/acs.jproteome.7b00370).  
30  
31 (21) Vargas, F., Weldon, K. C., Sikora, N., Wang, M., Zhang, Z., Gentry, E. C., Panitchpakdi, M. W., Caraballo,  
32 M., Dorrestein, P. C., Jarmusch, A. K. Protocol for Community-Created Public MS/MS Reference Library  
33 within the GNPS Infrastructure. *bioRxiv* **2019**, DOI: [10.1101/804401](https://doi.org/10.1101/804401).  
34  
35 (22) Grüning, B., Dale, R., Sjödin, A., Chapman, B. A., Rowe, J., Tomkins-Tinch, C. H., Valieris, R., Köster, J.,  
36 The Bioconda Team Bioconda: Sustainable and Comprehensive Software Distribution for the Life Sci-  
37 ences. *Nature Methods* **2018**, *15*, 475–476, DOI: [10.1038/s41592-018-0046-7](https://doi.org/10.1038/s41592-018-0046-7).  
38  
39 (23) RDKit: Open-source cheminformatics, <https://www.rdkit.org/>, (accessed: 2019-08-16).  
40  
41 (24) McKinney, W. In *Proceedings of the 9th Python in Science Conference*, ed. by van der Walt, S., Millman, J.,  
42 Austin, Texas, USA, 2010, pp 51–56.  
43  
44 (25) Pino, L., Lin, A., Bittremieux, W. 2018 YPIC Challenge: A Case Study in Characterizing an Unknown  
45 Protein Sample. *Journal of Proteome Research* **2019**, DOI: [10.1021/acs.jproteome.9b00384](https://doi.org/10.1021/acs.jproteome.9b00384).  
46  
47 (26) May, D. H., Bilmes, J., Noble, W. S. A Learned Embedding for Efficient Joint Analysis of Millions of Mass  
48 Spectra. *bioRxiv* **2018**, DOI: [10.1101/483263](https://doi.org/10.1101/483263).  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2 **For table of contents only**  
3  
4  
5



13  
14 **MS/MS spectrum processing & visualization**  
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