# 2018 YPIC Challenge: A case study in characterizing an unknown protein sample

**Lindsay Pino**[1], **Andy Lin**[1], **Wout Bittremieux**[*,1,2,3]

[1]Department of Genome Sciences, University of Washington, Seattle WA 98195, USA; [2]Department of Mathematics and Computer Science, University of Antwerp, 2020 Antwerp, Belgium; [3]Biomedical Informatics Network Antwerpen (biomina), 2020 Antwerp, Belgium

*Corresponding author: wout.bittremieux@uantwerpen.be, +32 3 265 34 07.

## Abstract

For the 2018 YPIC Challenge contestants were invited to try to decipher two unknown English questions encoded by a synthetic protein expressed in *Escherichia coli*. In addition to deciphering the sentence, contestants were asked to determine the 3D structure and determine any post-translation modifications left by the host organism.

We present how we analyzed this unknown sample using a tryptic digest with dynamic exclusion disabled to increase the signal-to-noise ratio of the measured molecules. Subsequently, spectral clustering was used to generate high-quality consensus spectra and condense the acquired MS/MS spectral data. *De novo* spectrum identification was used to determine the English questions encoded by the synthetic protein, and any post-translational modifications introduced by *E. coli* on the synthetic protein were detected using spectral networking.

Although the synthetic protein sample for the 2018 YPIC Challenge is not of biological interest, the experimental and computational strategy presented here can be directly used to analyze samples for which no protein sequence information is available or when the identity of the sample is unknown. All software and code to perform the bioinformatics analysis is available as open source, and a self-contained Jupyter notebook is provided to fully recreate the analysis.

## 1 Introduction

Mass spectrometry (MS) is a powerful analytical technique to characterize proteins in complex biological samples. The typical strategy to identify unknown tandem mass spectrometry (MS/MS) spectra is via sequence database searching [1]. Here, experimental MS/MS spectra are compared to theoretical spectra derived from a protein sequence database for the organism(s) of interest. Alternatively, spectral library searching can be used to identify unknown MS/MS spectra by comparing them against a library of high-quality, previously observed spectra with known peptide sequences [2, 3].

Both of these approaches depend on the availability of a ground truth reference set to which the unknown spectra are compared, either in the form of a sequence database or a spectral library. Alternatively, if such prior information is not available, such as, for example, during antibody sequencing or for non-model organisms whose genome

has not been sequenced yet, *de novo* searching can be used to directly derive peptide sequences from the unknown MS/MS spectra based on the mass differences between pairs of their fragment ion peaks [4].

Here, we describe our approach to characterize an unknown protein sample in the context of the 2018 Young Proteomics Investigators Club (YPIC) Challenge. YPIC is an initiative by the European Proteomics Association (EuPA) to connect and support young scientists in proteomics. As part of their activities they organize annual challenges where participants are invited to analyze mysterious protein samples [5].

The 2018 YPIC Challenge consisted of trying to decipher two unknown English questions encoded by a synthetic protein expressed in *E. coli*. The challenge encouraged participants to fully characterize the protein sample through several subtasks, such as protein sequence identification, detection of post-translational modifications (PTMs), and development of bioinformatic approaches.

Because the sample consisted of an unknown, synthetic, protein and no sequence database was available, we used *de novo* searching, in combination with spectral clustering, to identify the protein sequence. Additionally, spectral networking was used to discover common mass differences between spectra and detect potential PTMs. Finally, circular dichroism (CD) spectroscopy was used to analyze the protein's secondary structure.

All bioinformatics software that was used to analyze the data is freely available as open source. A self-contained Jupyter notebook [6] that contains all processing steps is available at `https://github.com/bittremieux/ypic_challenge_2018`, to fully reproduce the bioinformatics analysis.

# 2 Materials and methods

## 2.1 2018 YPIC Challenge description

We received a sample vial containing 12.5 µg of an unknown protein via mail from the organizers of the YPIC Challenge. As per the included product sheet, the synthetic protein was expressed in *E. coli* by PolyQuant and encoded two concatenated English questions. The sentence did not contain the letters 'B' and 'K', and the letters 'O' and 'U' were replaced by the letter 'K' in the protein. The protein sequence was flanked with 'MAGR' in the beginning and 'LAAALEHHHHHH' at the end for digestion and purification reasons.

The 2018 YPIC Challenge categories were as follows:

1. Answer *E. coli*'s question.

2. Three-dimensional grammar: Find out how this sentence folds.

3. Bioinformazing: Develop the coolest bioinformatics approach to decipher the sentence.

4. Protein punctuation: Look for the biological equivalent of punctuation: PTMs left behind by *E. coli*.

5. #Bioreactivity: Can you generate and describe bioreactivity in this Twitter-sized message?

## 2.2 Experimental procedures

### 2.2.1 Protein sample preparation

The sample was reconstituted with 125 µg 0.1 % formic acid (final concentration 0.1 µg/µL protein). An aliquot (1 µg; 10 µL) of reconstituted

sample was reduced (50 mM dithiothreitol), alkylated (150 mM iodoacetamide), and digested with Promega trypsin ($1 : 50$ enzyme—substrate ratio; 0.02 µg trypsin) for 4 h at 37 °C with shaking. Digested peptides were concentrated via speed-vac to a final concentration of 0.33 fmol/µL.

In addition to the conventional trypsin digest, following a CD spectroscopy solvent swap, the remaining sample was split into three parts and digested with three other proteases: pepsin, chymotrypsin, and Lys-C. The conditions for these reactions follow the trypsin digest conditions above, with the exception of the pepsin digestion which was held at a low pH (pH $< 2.0$).

### 2.2.2 LC-MS/MS data acquisition

Peptides were separated with a Waters NanoAcquity UPLC and emitted into a Thermo Q-Exactive HF tandem mass spectrometer. Pulled tip columns were created from 75 µm inner diameter fused silica capillary in-house using a laser pulling device and packed with 2.1 µm C18 beads (Dr. Maisch GmbH) to 300 mm. Trap columns were created from 150 µm inner diameter fused silica capillary fritted with Kasil on one end and packed with the same C18 beads to 25 mm. Buffer A was water and 0.1 % formic acid, while buffer B was 98 % acetonitrile and 0.1 % formic acid. For each injection, 3 µL of each sample was loaded with 5 µL 2 % B and eluted using the following program: 0 min to 90 min 2 % to 35 % B, 90 min to 100 min 35 % to 60 % B, followed by a 35 min washing gradient.

The Thermo Q-Exactive HF was set to positive mode in a top-20 configuration. Precursor scans ($300\,m/z$ to $2000\,m/z$) were collected at 60 000 resolution to hit an automatic gain control (AGC) target of $3 \times 10^6$. The maximum inject time was set to 100 ms. Fragment scans were collected at 30 000 resolution to hit an AGC target of $1 \times 10^5$ with a maximum inject time of 55 ms. The isolation width was set to 1.6 $m/z$ with a normalized collision energy of 27. Precursors with charge up to +6 that achieved a minimum AGC of $5 \times 10^3$ were acquired. Dynamic exclusion was disabled. The digested sample was acquired using this method in technical triplicate.

The mass spectrometry proteomics data have been deposited to the ProteomeXchange Consortium [7] via the PRIDE [8] partner repository with the dataset identifier PXD014003.

### 2.2.3 Circular dichroism spectroscopy

Following reconstitution of the protein sample as described above, the original protein sample,

minus the two μg of protein aliquoted for intact mass and trypsin digestion experiments, was speed vac'd to dryness to change to a CD spectroscopy buffer. The dried protein sample was reconstituted in 10 mM $KPO_4$ (pH 7.4) to 0.05 μg/μL (assuming 12.5 μg original protein per the product sheet and two μg used for the initial MS experiments) to meet the CD cuvette minimum volume requirement of 200 μL buffer.
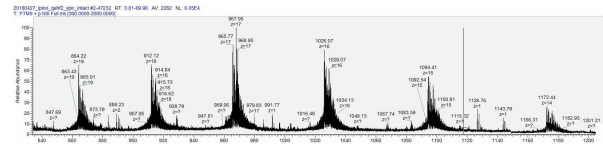
## 2.3 Data analysis

Raw files were converted to the MGF format using msconvert [9] for further processing. During conversion MS/MS spectra were centroided using the vendor algorithm and the precursor $m/z$ and charge was recalculated based on the preceding MS scan.

Next, MS/MS spectra were clustered and consensus spectra were generated using MaRaCluster [10] with a similarity p-value threshold of $10^{-5}$, precursor mass tolerance 50 ppm, and requiring at least 3 MS/MS spectra per cluster.

After spectral clustering low-quality clusters were removed by only retaining the clusters that represent at least 10 original spectra and whose consensus spectra have precursor charge 2 or 3.

The high-quality consensus spectra were used for *de novo* spectrum identification and spectral networking. DeNovoGUI [11] was used as a unified interface to the Novor [12], DirecTag [13], and PepNovo+ [14] *de novo* search engines. Settings for *de novo* spectrum identification were precursor mass tolerance 20 ppm; fragment mass tolerance 0.02 Da; and cysteine carbamidomethylation, methionine oxidation, and acetylation of the peptide N-terminus as variable modifications. Peptide–spectrum matches (PSMs) were visualized and manually investigated using DeNovoGUI.

A spectral network was constructed using the high-quality consensus spectra. Prior to matching spectra to each other they were preprocessed by removing noise peaks with an intensity below 5 % of the base peak intensity and at most the 150 most intense peaks were retained. Next, peak intensities were scaled by their square root before being normalized by their norm to have a magnitude of one. The shifted dot product [15] was used to match modified spectra to each other with fragment mass tolerance 0.02 Da. Each consensus spectrum formed a node in the spectral network, with an edge between two nodes if the shifted dot product between the two corresponding spectra was



**Figure 1:** MS1 scan of the intact synthetic protein indicating an approximate intact mass of 16.4 kDa.

greater than or equal to 0.8. Peptide sequences were assigned to nodes in the spectral network if the corresponding consensus spectra could be identified by Novor with a minimum score of 70. Only subgraphs in the spectral network consisting of at least three nodes were considered.

### 2.3.1 Code availability

Jupyter notebooks [6] containing all processing steps and analyses are available at https://github.com/bittremieux/ypic_challenge_2018. Custom processing is done in Python using open-source Python libraries including NumPy [16], pandas [17], NetworkX [18], Matplotlib [19], Seaborn [20], Pyteomics [21], and spectrum_utils [22]. The shifted dot product is implemented as an external C++ module for Python [15].
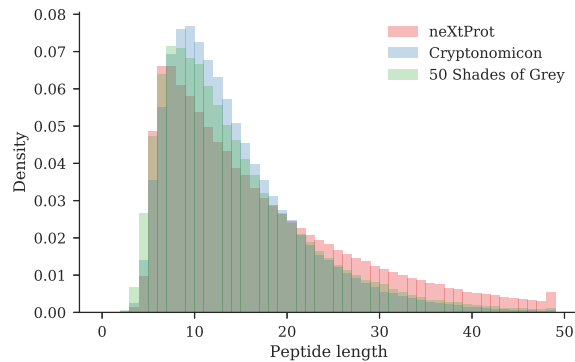
## 3 Results

## 3.1 Confirmation of intact mass

Prior to any peptide analysis, we determined the intact mass of the protein. While the final 2018 YPIC Challenge product sheet notes that the molecular weight of the protein is approximately 16.65 kDa, we received our challenge sample prior to the disclosure of this additional information. An MS1 spectrum of the intact mass confirms that the protein has an approximate mass of 16.4 kDa (figure 1).

## 3.2 Simulated digestion of English dictionary

When analyzing a protein of unknown sequence, one key decision is to determine which digestion enzyme to use. To help inform our decision we simulated the digestion of various corpuses using multiple proteases to determine whether they would generally yield peptides whose lengths are amenable to detection by mass spectrometry (figure 2).

**Figure 2:** Length of simulated tryptic peptides for various corpuses. NeXtProt is a database of human proteins, whereas Cryptonomicon and 50 Shades of Grey are two English fiction novels.

A simulated digestion of neXtProt [23], a database of human proteins, with trypsin while allowing for a single missed cleavage showed that a large portion of the resulting peptides will have a length between 6 and 10 amino acids, with the mode of the peptide length distribution at 7. While peptides of length 7 to 10 are perfectly amenable to detection by mass spectrometry, there is a significant tail in the distribution of peptide lengths. For example, approximately 15 % of the peptides consist of 30 or more amino acids, which is not ideal for detection by mass spectrometry.

An important issue with using neXtProt is that peptides in the human proteome are not expected to be a good proxy for 'peptides' found in human language. One possible reason is that the frequency of amino acids found in the human proteome is unlikely to be the same as the frequency of letters found in the English language. Therefore, we also simulated digestions of the text within two different English fiction novels with a varying vocabulary complexity: Cryptonomicon by Neal Stephenson [24] and 50 Shades of Grey by EL James [25]. Text files of the novels were found online and words containing the letters 'B' or 'K' were removed while the letters 'O' and 'U' were replaced by the letter 'K', in accordance with the challenge's instructions. Additionally, the letters 'J', 'X', and 'Z' were removed as these characters do not represent valid amino acids.

We found that a simulated digestion of these two novels yielded peptides whose lengths are slightly different than the length of peptides generated from neXtProt (figure 2). The majority of English peptides is slightly longer than those generated from neXtProt (the mode of the peptide

length distribution is 10 for Cryptonomicon and 7 for 50 Shades of Grey), while the English peptides include less very long peptides than neXtProt, as indicated by the right tail of the peptide length distributions. As a result, we found that trypsin is a suitable enzyme to digest the synthetic protein consisting of two English sentences. Additionally, we explored the digestion of neXtProt, Cryptonomicon, and 50 Shades of Grey with alternative proteases, including chymotrypsin, Glu-C, Lys-C, Arg-C, Asp-N, and pepsin, as well as combined digestions using two different proteases (supplementary **??**). These simulations again indicate that peptides generated from English are typically slightly longer than those generated from human proteins.

### 3.3 Synthetic protein identification

Since spectra were collected without dynamic exclusion enabled, molecules that are present in the sample will be selected multiple times for MS/MS measurement while spurious signals will only be measured a limited number of times. A downside of this approach is that the spectral data will contain multiple spectra that are virtually identical to each other as the same peptide is repeatedly measured. To condense the data volume the spectra were clustered with MaRaCluster [10]. Spectral clustering groups similar spectra together and creates a single consensus spectrum to represent each spectral cluster, reducing the number of spectra from 110 234 spectra in the original raw files to 380 consensus spectra representing at least ten spectra after spectral clustering (only retaining the spectra with precursor charge 2 or 3).

Next, these consensus spectra were identified. As no sequence database was available for the unknown synthetic protein *de novo* identification was performed. The Novor [12], DirecTag [13], and PepNovo+ [14] search engines were used through DeNovoGUI [11]. The resulting PSMs were subsequently manually validated, a task that became feasible thanks to the reduction in data volume by the spectral clustering. From the *de novo* identifications we were able to decode about 65 % of the unknown synthetic protein (based on the proportion of the mass of the identified peptides versus the mass of the intact protein). The following subsequences were compiled based on the identified PSMs, supplemented by some educated guesses:

- Start of protein: "Have you ever wondered what the mo[st]" (figures 3a to 3d)

- "[questio]ns in life ar[e]" (figure 3e)

- "[r]espect when it comes to what you" (figures 3f and 3g)

- End of protein: "[pro]duce in a cell." (figure 3h)

Unfortunately we were unable to identify the full synthetic protein as still a third of the sequence is missing. This is likely due to specific properties of the corresponding peptides which make them unamenable to identification using mass spectrometry, such as very short peptides after tryptic cleavage or peptides that cannot be properly ionized. We additionally tried to obtain complementary peptides using alternative proteases (pepsin, chymotrypsin, and Lys-C) to increase the sequence coverage. Unfortunately these experiments failed due to the sample loss observed during the preceding CD experiment (section 3.5).

## 3.4 Spectral networking to detect post-translational modifications

The typical approach to identify potentially modified peptides is by specifying variable modifications during a sequence database search. Similarly, variable modifications can be specified during *de novo* searching as well. However, *de novo* searching has to overcome several challenges compared to sequence database searching, including amino acid permutation complexity [4], and the inclusion of variable modifications exacerbates these challenges. Therefore, to maximize the confidence in the obtained *de novo* identifications only frequent PTMs introduced during sample processing [26] were specified to avoid a combinatorial explosion of the search space.

As an alternative strategy to find PTMs we have employed spectral networking [27]. A spectral network was constructed by representing each consensus spectrum as a node in a graph and connecting two nodes if their corresponding spectra are highly similar as measured by the shifted dot product [15, 28] (figure 4). Because the shifted dot product takes mass shifts induced by a modification into account while matching two spectra the spectral network will contain connections between modified peptides and their unmodified counterparts. Subsequently, based on the precursor mass difference between connected spectra in the spectral network and (partial) identifications of the spectra the presence and identity of various modifications, such as PTMs or amino acid substitutions, can be derived (figure 4).
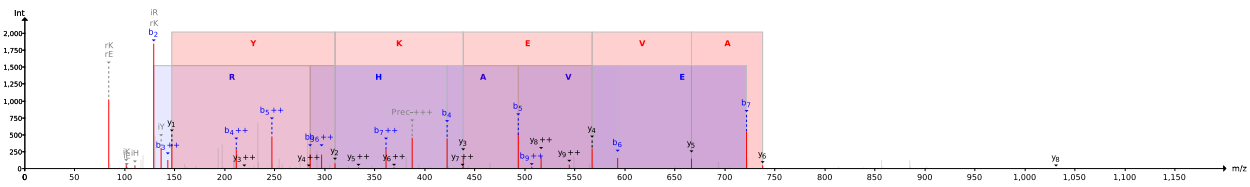
Connected spectra in the spectral network were manually checked for the presence of PTMs and the most frequently occurring mass differences were referenced to common modifications in Unimod [29]. This analysis indicated little to no systematic presence of PTMs. The most frequent mass differences were observed between unidentified spectra of low quality (manual quality assessment), likely derived from small molecular contaminants, and did not correspond to any common modifications. Although a more thorough investigation is recommended to conclusively determine the presence or absence of modifications, these preliminary results indicate that no PTMs are systematically introduced on the synthetic peptide by *E. coli*.

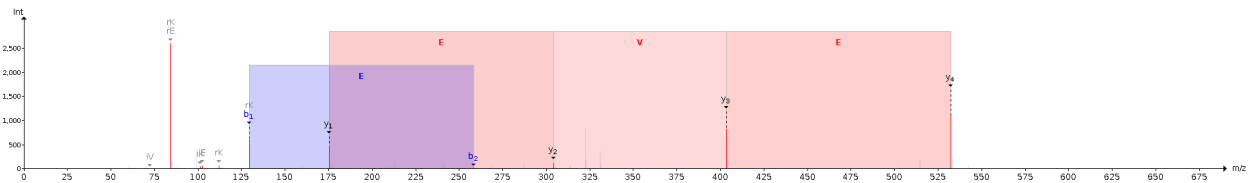## 3.5 Structural analysis using circular dichroism spectroscopy

We attempted CD spectroscopy to estimate the protein's secondary structure. The CD spectra, however, were inconclusive (data not shown). Based on absorption spectra acquired at the same time as the CD spectra, the concentration of protein in the CD cuvette was negligible. There are several reasons why the CD and absorption spectroscopy experiments might have failed. First, the concentration of protein ($0.05\,\mu g/\mu L$) may have been too dilute, considering the range of ideal protein concentration for CD spectroscopy is $0.1\,\mu g/\mu L$ to $0.2\,\mu g/\mu L$. Second, the buffer conditions used ($10\,\text{mM}\ KPO_4$ (pH 7.4)) may not be ideal for the protein's biochemistry, which would result in poor resolubilization of the protein. Third, the protein may have degraded during $-80\,°C$ storage and multiple freeze–thaw cycles during the course of the other experiments. Any one of these reasons may have contributed to the loss of protein observed in this experiment.
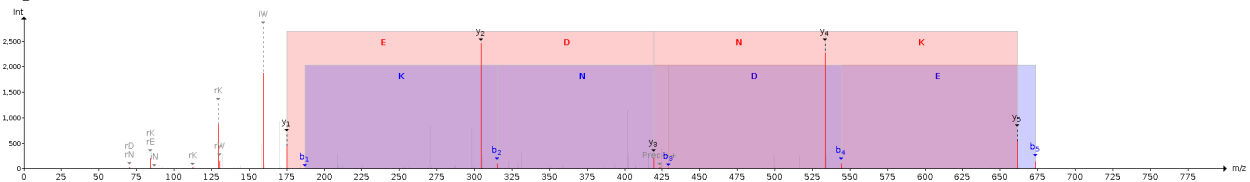
# 4 Conclusion

We have presented our results in identifying the unknown synthetic protein as part of the 2018 YPIC Challenge. Although we did not identify the full synthetic protein, based on a standard trypsin digest we are able to detect spectral evidence covering about two third of the unknown sequence. This is in line with the sequence coverage that is typically obtained during routine tryptic analyses of biological samples with a similar complexity. Although our attempts to use different pro-
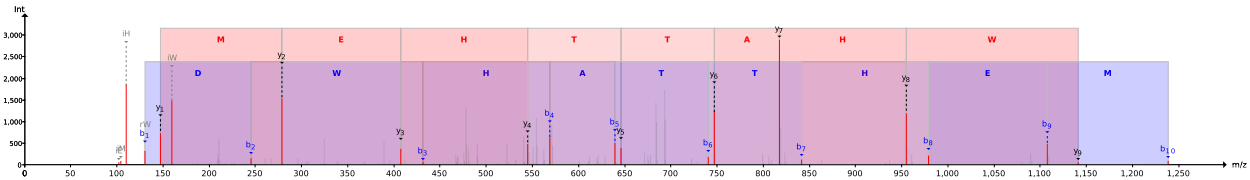
(**a**) Consensus spectrum 1945. Sequence: AGRHAVEYK, precursor mass: 386.88 $m/z$, precursor charge: 3, Novor score: 77.50, PepNovo+ score: 62.17.
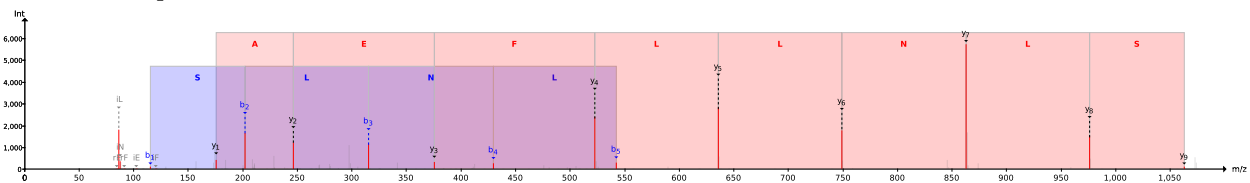


(**b**) Consensus spectrum 1503. Sequence: KEVER, precursor mass: 330.69 $m/z$, precursor charge: 2, Novor score: 92.20, PepNovo+ score: 70.47.
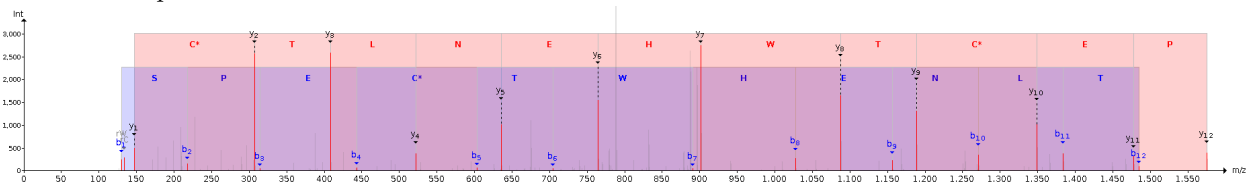


(**c**) Consensus spectrum 2136. Sequence: WKNDER, precursor mass: 424.21 $m/z$, precursor charge: 2, Novor score: 95.50, PepNovo+ score: 94.61.
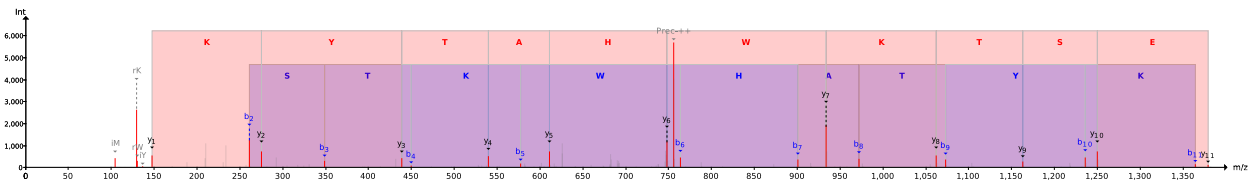


(**d**) Consensus spectrum 5178. Sequence: EDWHATTHEMK, precursor mass: 692.8 $m/z$, precursor charge: 2, Novor score: 88.70, PepNovo+ score: 139.17.
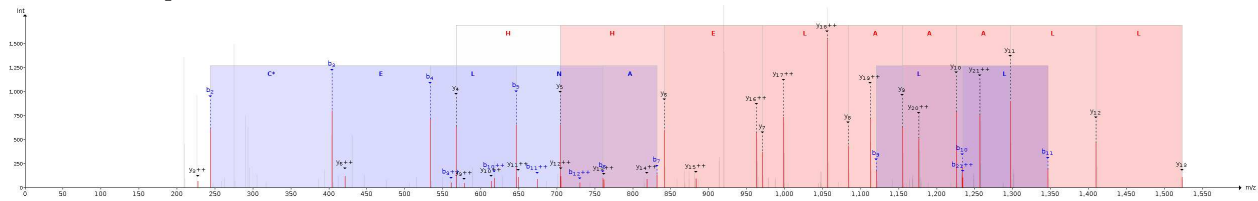


(**e**) Consensus spectrum 11694. Sequence: NSLNLLFEAR, precursor mass: 588.82 $m/z$, precursor charge: 2, Novor score: 94.60, PepNovo+ score: 122.66.



(**f**) Consensus spectrum 7109. Sequence: ESPECTWHENLTCK, precursor mass: 895.88 $m/z$, precursor charge: 2, Novor score: 94.10, PepNovo+ score: 192.51.

(**g**) Consensus spectrum 7109. Sequence: MESTKWHATYKK, precursor mass: $755.38\,m/z$, precursor charge: 2, Novor score: 93.20, PepNovo+ score: 156.91.



(**h**) Consensus spectrum 9658. Sequence: DKCELNACELLLAAALEHHDYNR, precursor mass: $919.1\,m/z$, precursor charge: 3, Novor score: 57.10.

**Figure 3:** Relevant PSMs decoding the unknown synthetic protein.

teases to increase the sequence coverage failed due to lack of sample material and sample loss that occurred during multiple experiments, we anticipate that this strategy would have generated alternative peptides [30]. Additionally, using unconventional digestion strategies such as microwave-assisted digestion to obtain semi-random peptide cleavage [31], might have increased the protein sequence coverage.

Despite lacking the full protein sequence, we used spectral clustering and spectral networking to investigate the presence of frequent modifications. Based on this analysis we did not see any systematic modifications on the synthetic protein. This corresponds to the lack of notable PTMs in *E. coli* as well, although a more extensive analysis is recommended to conclusively determine the absence of any modifications.

Although in this case the sample consisted of a contrived synthetic protein in the context of the 2018 YPIC Challenge, the experimental and computational strategy we have described here can similarly be used to analyze other unknown protein samples that are of more biological interest, such as, for example, antibody sequencing. Notably, our spectral clustering approach can be used to increase the signal-to-noise ratio of spectra prior to *de novo* identification [32]. Additionally, spectral networking is an increasingly popular strategy to analyze small molecules measured by mass spectrometry [33].
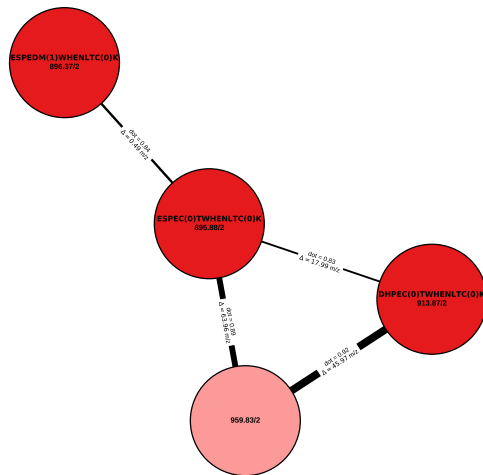
## Acknowledgement

## References

[1] Eng, J. K. et al. "A Face in the Crowd: Recognizing Peptides through Database Search." In: *Molecular & Cellular Proteomics* 10.11 (Nov. 1, 2011), R111.009522. DOI: `10.1074/mcp.R111.009522`.

[2] Griss, J. "Spectral Library Searching in Proteomics." In: *PROTEOMICS* 16.5 (Mar. 2016), pp. 729–740. DOI: `10.1002/pmic.201500296`.

[3] Shao, W. and Lam, H. "Tandem Mass Spectral Libraries of Peptides and Their Roles in Proteomics Research." In: *Mass Spectrometry Reviews* 36.5 (Sept. 2017), pp. 634–648. DOI: `10.1002/mas.21512`.

[4] Muth, T. et al. "A Potential Golden Age to Come-Current Tools, Recent Use Cases, and Future Avenues for de Novo Sequencing in Proteomics." In: *PROTEOMICS* 18.18 (Sept.

**Figure 4:** A spectral network connects (un)modified spectra. The peptide sequence (if known) and the precursor mass and precursor charge are shown for each node in the spectral network. Edges between two nodes are annotated with the corresponding precursor mass difference. The spectral similarity based on the shifted dot product is indicated by the weight of the edge.

The spectral network shows a strong similarity between multiple spectra despite small differences in the identified sequences due to amino acid substitutions (CT ↔ DM, ES ↔ DH). Although the spectrum corresponding to the light red shaded node could not be fully identified through *de novo* searching, its high similarity to related spectra indicates that it was likely derived from the same peptide. Indeed, a full identification was precluded by the absence of any successfully matched b-ions, while the C-terminal tag "CTWHENLTCK" could still be annotated based on the y-ions.

2018), p. 1700150. DOI: 10 . 1002 / pmic . 201700150.

[5] Indeykina, M. I., Podgrudkov, D. A., and Kononikhin, A. S. "The Author Identified by His Method: EuPA YPIC Challenge Solved." In: *EuPA Open Proteomics* 20 (Dec. 2018), pp. 1–8. DOI: 10.1016/j.euprot.2018.10.001.

[6] Thomas, K. et al. "Jupyter Notebooks – A Publishing Format for Reproducible Computational Workflows." In: *Positioning and Power in Academic Publishing: Players, Agents and Agendas*. IOS Press, 2016, pp. 87–90.

[7] Deutsch, E. W. et al. "The ProteomeXchange Consortium in 2017: Supporting the Cultural Change in Proteomics Public Data Deposition." In: *Nucleic Acids Research* 45.D1 (Jan. 4, 2017), pp. D1100–D1106. DOI: 10 . 1093/nar/gkw936.

[8] Perez-Riverol, Y. et al. "The PRIDE Database and Related Tools and Resources in 2019: Improving Support for Quantification Data." In: *Nucleic Acids Research* 47.D1 (Jan. 8, 2019), pp. D442–D450. DOI: 10.1093/nar/gky1106.

[9] Chambers, M. C. et al. "A Cross-Platform Toolkit for Mass Spectrometry and Proteomics." In: *Nature Biotechnology* 30.10 (Oct. 10, 2012), pp. 918–920. DOI: 10.1038/nbt.2377.

[10] The, M. and Käll, L. "MaRaCluster: A Fragment Rarity Metric for Clustering Fragment Spectra in Shotgun Proteomics." In: *Journal of Proteome Research* 15.3 (Mar. 4, 2016), pp. 713–720. DOI: 10.1021/acs.jproteome.5b00749.

[11] Muth, T. et al. "DeNovoGUI: An Open Source Graphical User Interface for *de Novo* Sequencing of Tandem Mass Spectra." In: *Journal of Proteome Research* 13.2 (Feb. 7, 2014), pp. 1143–1146. DOI: 10 . 1021 / pr4008078.

[12] Ma, B. "Novor: Real-Time Peptide de Novo Sequencing Software." In: *Journal of The American Society for Mass Spectrometry* 26.11 (Nov. 2015), pp. 1885–1894. DOI: 10.1007/s13361-015-1204-0.

[13] Tabb, D. L. et al. "DirecTag: Accurate Sequence Tags from Peptide MS/MS through Statistical Scoring." In: *Journal of Proteome Research* 7.9 (Sept. 5, 2008), pp. 3838–3846. DOI: 10.1021/pr800154p.

[14] Frank, A. et al. "Peptide Sequence Tags for Fast Database Search in Mass-Spectrometry." In: *Journal of Proteome Research* 4.4 (Aug. 8, 2005), pp. 1287–1295. DOI: `10.1021/pr050011x`.

[15] Bittremieux, W. et al. "Fast Open Modification Spectral Library Searching through Approximate Nearest Neighbor Indexing." In: *Journal of Proteome Research* 17.10 (Oct. 5, 2018), pp. 3463–3474. DOI: `10.1021/acs.jproteome.8b00359`.

[16] Van der Walt, S., Colbert, S. C., and Varoquaux, G. "The NumPy Array: A Structure for Efficient Numerical Computation." In: *Computing in Science & Engineering* 13.2 (Mar. 2011), pp. 22–30. DOI: `10.1109/MCSE.2011.37`.

[17] McKinney, W. "Data Structures for Statistical Computing in Python." In: *Proceedings of the 9th Python in Science Conference*. Ed. by van der Walt, S. and Millman, J. Austin, Texas, USA, 2010, pp. 51–56.

[18] Hagberg, A. A., Schult, D. A., and Swart, P. J. "Exploring Network Structure, Dynamics, and Function Using NetworkX." In: *Proceedings of the 7th Python in Science Conference - SciPy '08*. Ed. by Varoquaux, G., Vaught, T., and Millman, J. Pasadena, CA USA, 2008, pp. 11–15.

[19] Hunter, J. D. "Matplotlib: A 2D Graphics Environment." In: *Computing in Science & Engineering* 9.3 (June 18, 2007), pp. 90–95. DOI: `10.1109/MCSE.2007.55`.

[20] Waskom, M. et al. *Mwaskom/Seaborn: V0.8.1 (September 2017)*. Sept. 3, 2017. DOI: `10.5281/zenodo.883859`.

[21] Goloborodko, A. A. et al. "Pyteomics-a Python Framework for Exploratory Data Analysis and Rapid Software Prototyping in Proteomics." In: *Journal of The American Society for Mass Spectrometry* 24.2 (Feb. 1, 2013), pp. 301–304. DOI: `10.1007/s13361-012-0516-6`.

[22] *bittremieux/spectrum_utils: Simple MS/MS spectrum preprocessing and visualization in Python*. `https://github.com/bittremieux/spectrum_utils`. (accessed: 2019-02-22).

[23] Gaudet, P. et al. "The neXtProt Knowledgebase on Human Proteins: Current Status." In: *Nucleic Acids Research* 43.D1 (Jan. 28, 2015), pp. D764–D770. DOI: `10.1093/nar/gku1178`.

[24] Stephenson, N. *Cryptonomicon*. 1st ed. New York: Avon Press, 1999. 918 pp.

[25] James, E. L. *Fifty Shades of Grey*. Fifty shades trilogy 1. New York: Vintage Books, 2015. 514 pp.

[26] Bittremieux, W. et al. "Quality Control in Mass Spectrometry-Based Proteomics." In: *Mass Spectrometry Reviews* 37.5 (Sept. 2018), pp. 697–711. DOI: `10.1002/mas.21544`.

[27] Bandeira, N. et al. "Protein Identification by Spectral Networks Analysis." In: *Proceedings of the National Academy of Sciences* 104.15 (Apr. 10, 2007), pp. 6140–6145. DOI: `10.1073/pnas.0701130104`.

[28] Bittremieux, W., Laukens, K., and Noble, W. S. "Extremely Fast and Accurate Open Modification Spectral Library Searching of High-Resolution Mass Spectra Using Feature Hashing and Graphics Processing Units." In: *bioRxiv* (May 5, 2019). DOI: `10.1101/627497`.

[29] Creasy, D. M. and Cottrell, J. S. "Unimod: Protein Modifications for Mass Spectrometry." In: *PROTEOMICS* 4.6 (Apr. 5, 2004), pp. 1534–1536. DOI: `10.1002/pmic.200300744`.

[30] Tsiatsiani, L. and Heck, A. J. R. "Proteomics beyond Trypsin." In: *FEBS Journal* 282.14 (July 2015), pp. 2612–2626. DOI: `10.1111/febs.13287`.

[31] Savidor, A. et al. "Database-Independent Protein Sequencing (DiPS) Enables Full-Length de Novo Protein and Antibody Sequence Determination." In: *Molecular & Cellular Proteomics* 16.6 (June 1, 2017), pp. 1151–1161. DOI: `10.1074/mcp.O116.065417`.

[32] Griss, J. et al. "Spectral Clustering Improves Label-Free Quantification of Low-Abundant Proteins." In: *Journal of Proteome Research* 18.4 (Apr. 5, 2019), pp. 1477–1485. DOI: `10.1021/acs.jproteome.8b00377`.

[33] Wang, M. et al. "Sharing and Community Curation of Mass Spectrometry Data with Global Natural Products Social Molecular Networking." In: *Nature Biotechnology* 34.8 (Aug. 9, 2016), pp. 828–837. DOI: `10.1038/nbt.3597`.