

¹ Peptacular: A Python Library for ProForma 2.0-Compliant Peptide Analysis

³ Patrick Tyler Garrett  ¹ and John R. Yates III  

⁴ 1 The Scripps Research Institute, United States  Corresponding author

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Open Journals](#) 

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

⁵ Summary

⁶ Mass spectrometry-based proteomics relies on computational methods to identify and ⁷ characterize amino acid sequences (Angel et al., 2012). These sequences, which range ⁸ from short peptides to complete proteins, exhibiting considerable chemical complexity. This ⁹ complexity includes post-translational modifications, variable charge states, charge adducts, ¹⁰ neutral losses, and isotopic patterns (Smith & Kelleher, 2013). ProForma notation provides ¹¹ a standardized, human-readable representation of this complexity, facilitating consistent ¹² communication of peptide and protein identifications across the proteomics community (LeDuc ¹³ et al., 2018).

¹⁴ Peptacular is a pure Python library designed for working with ProForma notation in ¹⁵ high-throughput proteomics applications. The library provides functionality for sequence ¹⁶ manipulation, molecular mass and composition calculation, physicochemical property prediction, ¹⁷ fragment ion generation, in silico enzymatic digestion, isotopic distribution computation, and ¹⁸ format conversion between ProForma and vendor-specific notations.

¹⁹ Statement of Need

²⁰ The proteomics field has historically lacked standardization for peptide and protein sequences. ²¹ Individual software tools have implemented proprietary notations, creating barriers to data ²² integration and reanalysis. While the Proteomics Standards Initiative developed ProForma ²³ notation to address this issue, adoption remains limited due to its relative novelty and the ²⁴ scarcity of computational tools supporting the standard.

²⁵ Modern proteomics experiments routinely identify hundreds of thousands of peptide-spectrum ²⁶ matches. Results are typically exported as tabular files (CSV, TSV, Parquet) compatible with ²⁷ Python data analysis libraries such as **pandas** (Team, 2025) and **polars** (Vink et al., 2025). ²⁸ However, existing proteomics libraries lack efficient mechanisms for performing sequence-based ²⁹ calculations at this scale or integrating seamlessly with dataframe-centric workflows.

³⁰ Peptacular addresses these limitations by providing two complementary APIs: an object-oriented ³¹ interface, and a functional interface. The functional API accepts serialized ProForma sequences, ³² applies transformations using parallel processing, and returns serialized results. This design ³³ enables direct integration with dataframes without sacrificing computational performance. The ³⁴ object-oriented interface allows for direct manipulation of the annotation objects without the ³⁵ need for serialization or parsing, improving performance for more complex operations.

³⁶ Features and Implementation

³⁷ Peptacular's employs lazy evaluation to minimize memory overhead: modifications remain in ³⁸ serialized form until explicitly required for calculations. Since Python strings are immutable

39 and proteomics datasets typically contain repeated modifications, Peptacular utilizes caching
40 wherever possible, yielding performance and memory improvements.

41 The object-oriented API uses a factory pattern to construct Annotation objects, enabling
42 sequence manipulation including modification placement, isotopic labeling, and custom fragment
43 ion definitions. The functional API provides stateless functions that operate directly on serialized
44 sequences, facilitating parallel execution across CPU cores.

45 Modification reference data, including masses, compositions, and identifiers from Unimod
46 ([Creasy & Cottrell, 2004](#)), and other databases, are embedded within the package as Python
47 modules rather than external files. This design eliminates file I/O overhead for parsing
48 the supported ontologies, improving performance. Additionally, Peptacular uses conditional
49 initialization: modification-related data structures are only instantiated when the corresponding
50 modification type is present in the input sequences, reducing startup overhead and memory
51 consumption.

52 The library supports bidirectional conversion between ProForma notation and common vendor-
53 specific formats, facilitating integration with existing proteomics workflows while promoting
54 migration toward standardized representations.

55 Related Software

56 Several existing tools address aspects of peptide manipulation and property calculation, including
57 pyteomics ([Goloborodko et al., 2013](#)) and biopython ([Cock et al., 2009](#)). However, these tools
58 were developed prior to the establishment of ProForma 2.0 notation and consequently lack
59 native support for many modern ProForma-specific features. Furthermore, these packages
60 include extensive dependencies and functionality that extends well beyond the scope of peptide
61 notation standardization. Peptacular addresses this gap by providing a focused implementation
62 designed to serve as computational infrastructure connecting different software tools and
63 packages within the proteomics ecosystem.

64 Proforma Compliance

65 5.3.1 Base-ProForma Compliance

- 66 **Amino acids (+UO)** - AAHCFKUOT (§6.1)
- 67 **Unimod names** - PEM[Oxidation]AT (§6.2.1)
- 68 **PSI-MOD names** - PEM[monohydroxylated residue]AT (§6.2.1)
- 69 **Unimod numbers** - PEM[UNIMOD:35]AT (§6.2.2)
- 70 **PSI-MOD numbers** - PEM[MOD:00425]AT (§6.2.2)
- 71 **Delta masses** - PEM[+15.995]AT (§6.2.3)
- 72 **N-terminal modifications** - [Carbamyl]-QPEPTIDE (§6.3)
- 73 **C-terminal modifications** - PEPTIDEG-[Methyl] (§6.3)
- 74 **Labile modifications** - {Glycan:Hex}EM[U:Oxidation]EV (§6.4)
- 75 **Multiple modifications** - MPGNW[Oxidation][Carboxymethyl]PESQE (§6.5)
- 76 **Information tag** - ELV[INFO:AnyString]IS (§6.6)

77 5.3.2 Level 2-ProForma Compliance

- 78 **Ambiguous amino acids** - BZJX (§7.1)
- 79 **Prefixed delta masses** - PEM[U:+15.995]AT (§7.2)
- 80 **Mass gap** - PEX[+147.035]AT (§7.3)
- 81 **Formulas** - PEM[Formula:0]AT, PEM[Formula:[1701]]AT (§7.4)
- 82 **Mass with interpretation** - PEM[+15.995|Oxidation]AT (§7.5)
- 83 **Unknown mod position** - [Oxidation]?PEMAT (§7.6.1)

- 84 **Set of positions** - PEP[Oxidation#1]M[#1]AT (§7.6.2)
- 85 **Range of positions** - PRT(ESFRMS)[+19.0523]ISK (§7.6.3)
- 86 **Position scores** - PEP[Oxidation#1(0.95)]M[#1(0.05)]AT (§7.6.4)
- 87 **Range position scores** - (PEP)[Oxidation#1(0.95)]M[#1(0.05)]AT (§7.6.5)
- 88 **Amino acid ambiguity** - (?VCH)AT (§7.7)
- 89 **Modification prefixes** - PEPM[U:Oxidation]AS[M:0-phospho-L-serine] (§7.8)

90 5.3.3 Level 2-ProForma + Top-Down

- 91 **RESID modifications** - EM[R:L-methionine sulfone]EM[RESID:AA0581] (§8.1)
- 92 **Names** - (>Heavy chain)EVQLVESG (§8.2)

93 5.3.4 Level 2-ProForma + Cross-Linking

- 94 **XL-MOD modifications** - EVTK[X:Aryl azide]LEK[XLMOD:00114]SEFD (§9.1)
- 95 **Cross-linkers (intrachain)** - EVTK[X:Aryl azide#XL1]LEK[#XL1]SEFD (§9.2.1)
- 96 **Cross-linkers (interchain)** - EVTK[X:Aryl azide#XL1]L//EK[#XL1]SEFD (§9.2.2)
- 97 **Branches** - ED[MOD:00093#BRANCH]//D[#BRANCH]ATR (§9.3)

98 5.3.5 Level 2-ProForma + Glycans

- 99 **GNO modifications** - NEEYN[GNO:G59626AS]K (§10.1)
- 100 **Glycan compositions** - NEEYN[Glycan:Hex5HexNAc4NeuAc1]K (§10.2)

101 5.3.6 Level 2-ProForma + Advanced Complexity

- 102 **Charged formulas** - SEQUEN[Formula:Zn1:z+2]CE (§11.1)
- 103 **Controlling placement** - PTI(MERMERME)[+32|Position:E]PTIDE (§11.2)
- 104 **Global isotope** - <13C>CARBON (§11.3.1)
- 105 **Fixed modifications** - <[Oxidation]@M>ATPEMILTCMGCLK (§11.3.2)
- 106 **Chimeric spectra** - NEEYN+SEQUEN (§11.4)
- 107 **Charges** - SEQUEN/2, SEQUEN/[Na:z+1,H:z+1] (§11.5)
- 108 **Ion notation** - SEQUEN-[b-type-ion] (§11.6)

109 AI usage disclosure

110 Various AI models were used to aid in the creation of this softwares codebase, tests, and
 111 documentation, mainly Claude Sonnet 4.5, Gemini 3 Pro, and Github's copilot autocomplete
 112 extension. The AI models were used via the copilot extension in VS code. No AI models were
 113 used in the create of this manuscript.

114 Angel, T. E., Aryal, U. K., Hengel, S. M., Baker, E. S., Kelly, R. T., Robinson, E. W., &
 115 Smith, R. D. (2012). Mass spectrometry-based proteomics: existing capabilities and future
 116 directions. *Chemical Society Reviews*, 41(10), 3912. <https://doi.org/10.1039/c2cs15331a>

117 Cock, P. J. A., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., Friedberg,
 118 I., Hamelryck, T., Kauff, F., Wilczynski, B., & De Hoon, M. J. L. (2009). Biopython:
 119 freely available Python tools for computational molecular biology and bioinformatics.
Bioinformatics, 25(11), 1422–1423. <https://doi.org/10.1093/bioinformatics/btp163>

121 Creasy, D. M., & Cottrell, J. S. (2004). Unimod: Protein modifications for mass spectrometry.
 122 *PROTEOMICS*, 4(6), 1534–1536. <https://doi.org/10.1002/pmic.200300744>

123 Goloborodko, A. A., Levitsky, L. I., Ivanov, M. V., & Gorshkov, M. V. (2013). Pyteomics—A
 124 Python framework for exploratory data analysis and rapid software prototyping in proteomics.
 125 *Journal of the American Society for Mass Spectrometry*, 24(2), 301–304. <https://doi.org/10.1007/s13361-012-0516-6>

- 127 LeDuc, R. D., Schwämmle, V., Shortreed, M. R., Cesnik, A. J., Solntsev, S. K., Shaw, J.
128 B., Martin, M. J., Vizcaino, J. A., Alpi, E., Danis, P., Kelleher, N. L., Smith, L. M.,
129 Ge, Y., Agar, J. N., Chamot-Rooke, J., Loo, J. A., Pasa-Tolic, L., & Tsybin, Y. O.
130 (2018). ProForma: a standard proteoform notation. *Journal of Proteome Research*, 17(3),
131 1321–1325. <https://doi.org/10.1021/acs.jproteome.7b00851>
- 132 Smith, L. M., & Kelleher, N. L. (2013). Proteoform: a single term describing protein complexity.
133 *Nature Methods*, 10(3), 186–187. <https://doi.org/10.1038/nmeth.2369>
- 134 Team, P. D. (2025). pandas-dev/pandas: Pandas. Zenodo (CERN European Organization for
135 Nuclear Research). <https://doi.org/10.5281/zenodo.17992932>
- 136 Vink, R., De Gooijer, S., Beedie, A., Burghoorn, G., Nameexhaustion, Peters, O., Gorelli,
137 M. E., Reswqa, Van Zundert, J., Marshall, Hulselmans, G., Grinstead, C., Denecker, K.,
138 Manley, L., chielP, Turner-Trauring, I., Valtar, K., Mitchell, L., Sprenkels, A., ... Magarick,
139 J. (2025). pola-rs/polars: Python Polars 1.36.1. Zenodo (CERN European Organization
140 for Nuclear Research). <https://doi.org/10.5281/zenodo.17873635>