

teiphy: General-purpose Python utility for converting TEI XML collations to NEXUS and other formats

Joey McCollum¹ and Robert Turnbull²

¹ Institute for Religion and Critical Inquiry, Australian Catholic University, Australia ² Melbourne Data Analytics Platform, University of Melbourne, Australia

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Published: unpublished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](#)).

Summary

Textual scholars have been using phylogenetics to analyze manuscript traditions since the early 1990s ([Robinson & O'Hara, 1992](#)). Many standard phylogenetic software packages accept as input the NEXUS file format ([Maddison et al., 1997](#)). The teiphy program takes a collation of texts encoded using the Text Encoding Initiative (TEI) guidelines and converts it to a NEXUS format so that it can be used for phylogenetic analysis. It can convert to other formats, as well.

Statement of Need

The TEI aims to provide an international standard for digital encoding textual information for the humanities ([Ide & Sperberg-McQueen, 1995](#)). The TEI guidelines describe an XML format for encoding a critical apparatus ([Consortium, 2022](#)). Due to its rich and well-documented set of elements for expressing a wide range of features in manuscript transcriptions, collations, and critical editions, TEI XML has become the *de facto* format for textual data in the digital humanities ([Fischer, 2020](#)). Its expressive power has proven increasingly valuable since its release, as scholars have learned—sometimes the hard way—that digitized texts should (1) preserve as much detail as they can from their material sources, including paratextual features; (2) reproduce the text of their sources as closely as possible, with editorial regularizations to things like orthography, accentuation, and scribal shorthand encoded alongside rather than in place of the source text; and (3) uncertainties about what a source read, both in transcriptions and collations, should be described as accurately as possible, allowing for degrees of uncertainty and multiple choices for disambiguations if necessary.

Such principles have much bearing on the editing of critical texts, a task fundamental to both digital humanities and classical philology. Within the digital humanities, phylogenetic algorithms have been popular approaches to this task. Taking the most arduous part of reconstructing a textual tradition and delegating it to a computer proved to be a promising technique, and its successful demonstration with a portion of *The Canterbury Tales* was a milestone in the development of the field ([Barbrook et al., 1998](#)). Soon after this, the same methods were applied more comprehensively to the tradition of *Lanseloet van Denemerken* in a work that would formalize many practical rules for computer-assisted textual criticism ([Salemans, 2000](#)). Since then, phylogenetic methods have quickly evolved ([Felsenstein, 2004](#)), and textual critics have adapted the improvements and even added their own innovations to make the process more suitable for their purposes ([Carlson, 2015](#); [Edmondson, 2019](#); [Hyytiäinen, 2021](#); [Spencer et al., 2002, 2004](#); [Turnbull, 2020](#)).

As their name might suggest, phylogenetic methods originated in the setting of evolutionary biology. They have a natural place in textual criticism given the deep analogy between the two fields: a sequence alignment, which consists of taxa, sites or characters, and the states of taxa

at those characters, corresponds almost identically to a collation, which consists of witnesses to the text, locations of textual variation (which we will call “variation units” from here on), and the variant readings attested by witnesses at those points.

Most phylogenetic software, however, expects inputs not in TEI XML format, but in NEXUS format (Maddison et al., 1997). This format was conceived with versatility in mind, and this design choice has been vindicated in its applicability with textual data, but NEXUS is not equipped or meant to express the same kinds of details that TEI XML is. Conversely, for those interested primarily in working with the collation as an alignment, TEI XML is overkill. Thus, a great chasm has been fixed between the two formats, and the only way to cross over it is by conversion.

Another format of value for text-critical phylogenetics is the input format associated with the STEMMMA software developed by Stephen C. Carlson for his 2012 thesis (Carlson, 2015); the code is hosted at <https://github.com/stemmatic/stemmma>. Carlson’s software expands on traditional maximum parsimony-based phylogenetic algorithms with rules to account for contamination or mixture in the manuscript tradition. While it has so far only been applied to books of the New Testament, it is just as applicable to other traditions, and a way of converting TEI XML collations of other texts to a format that can be used by this software could help bridge this gap.

Finally, it is also worth mentioning that other basic machine-learning approaches to textual criticism, which are frequently based on clustering and biclustering algorithms (Finney, 2018; McCollum, 2019; [thorpe_multivariate2002?](#)), expect the collation data to be encoded as a matrix with a row for each variant reading and a column for each witness. Thus, a means of converting the essential data from TEI XML collation to a NumPy array (Harris et al., 2020) and other related formats is a need for applications like these.

Design

While the conversion process is a straightforward one for most collation data, lacunae, retroversions, and other sources of ambiguity occasionally make a one-to-one mapping of witnesses to readings impossible, and in some cases, one disambiguation may be more likely than another in a quantifiable way. Mechanisms for accommodating such situations exist in both TEI XML and NEXUS, and for likelihood-based phylogenetic methods, “soft decisions” about the states at the leaves and even the root of the tree can provide useful information to the inference process. For these reasons, we wanted to ensure that these types of judgments, as well as other rich features from TEI XML, could be respected (and, where, necessary, preserved) in the conversion process.

Collations should preserve as much detail as possible, including information on how certain types of data can be normalized and collapsed for analysis. Since one might want to conduct the same analysis at different levels of granularity, the underlying collation data should be available for use in any case, and only the output should reflect changes in the desired level of detail. Likewise, as noted in the previous section, uncertainty about witnesses’ attestations should be encoded in the collation and preserved in the conversion of the collation.

For text-critical purposes, differences in granularity typically concern which types of variant readings we consider important for analysis. At the lowest level, readings with uncertain or reconstructed portions are almost always considered identical with their reconstructions (provided these reconstructions can be made unambiguously) for the purpose of analysis. Defective forms that are obvious misspellings of a more substantive reading are often treated the same way. Even orthographic subvariants that reflect equally “correct” regional spelling practices may be considered too common and of too trivial a nature to be of value for analysis. Other readings that do not fall under these rubrics but are nevertheless considered manifestly secondary (due to late and/or isolated attestation, for instance), may also be considered uninformative “noise” that is better left filtered out.

Use Case

Ephesians UBS example? Is there an example from

Availability

The software can be installed through the Python Package Index (PyPI) and the source code is available under the MIT license from the GitHub repository. The automated testing suite has 100% coverage.

Acknowledgements

The authors wish to thank Stephen C. Carlson for his feedback on this project.

Funding

This work was supported by an Australian Government Research Training Program (RTP) Scholarship.

References

- Barbrook, A. C., Howe, C. J., Blake, N., & Robinson, P. (1998). The Phylogeny of The Canterbury Tales. *Nature*, 394, 839. <https://doi.org/10.1038/29667>
- Carlson, S. C. (2015). *The Text of Galatians and Its History*. Mohr Siebeck. ISBN: 978-3-16-153323-5
- Consortium, T. (2022). *TEI P5: Guidelines for Electronic Text Encoding and Interchange: Critical Apparatus [v.4.4.0]*. <https://www.tei-c.org/release/doc/tei-p5-doc/en/html/TC.html>.
- Edmondson, A. C. (2019). *An Analysis of the Coherence-Based Genealogical Method Using Phylogenetics*. University of Birmingham. <http://etheses.bham.ac.uk/id/eprint/9150>
- Felsenstein, J. (2004). *Inferring Phylogenies*. Sinauer Associates.
- Finney, T. J. (2018). How to Discover Textual Groups. *Digital Studies/Le Champ Numérique*, 8. <https://doi.org/10.16995/dscn.291>
- Fischer, F. (2020). Representing the Critical Text. In P. Roelli (Ed.), *Handbook of stemmatology: History, methodology, digital approaches* (pp. 405–427). De Gruyter.
- Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., Kerkwijk, M. H. van, Brett, M., Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., ... Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, 585, 357–362. <https://doi.org/10.1038/s41586-020-2649-2>
- Hyytiäinen, P. (2021). The Changing Text of Acts: A Phylogenetic Approach. *TC: A Journal of Biblical Textual Criticism*, 26, 1–28.
- Ide, N. M., & Sperberg-McQueen, C. M. (1995). The TEI: History, Goals and Future. *Computers and the Humanities*, 29(1), 5–15.
- Maddison, D. R., Swofford, D. L., & Maddison, W. P. (1997). NEXUS: An Extensible File Format for Systematic Information. *Systematic Biology*, 46(4), 590–621. <https://doi.org/10.1093/sysbio/46.4.590>

- 130 McCollum, J. (2019). Biclustering Readings and Manuscripts via Non-negative Matrix Factor-
131 ization, with Application to the Text of Jude. *Andrews University Seminary Studies*, 57(1),
132 61–89.
- 133 Robinson, P., & O'Hara, R. J. (1992). Report on the Textual Criticism Challenge 1991. *Bryn*
134 *Mawr Classical Review*, 3(4), 331–337.
- 135 Salemans, B. J. P. (2000). *Building Stemmas with the Computer in a Cladistic, Neo-*
136 *Lachmannian, Way: The Case of Fourteen Text Versions of Lanceloet van Denemerken*.
137 Katholieke Universiteit Nijmegen. <https://hdl.handle.net/2066/147058>
- 138 Spencer, M., Wachtel, K., & Howe, C. J. (2002). The Greek Vorlage of the Syra Harclensis: A
139 Comparative Study on Method in Exploring Textual Genealogy. *TC: A Journal of Biblical*
140 *Textual Criticism*, 7. <http://jbtc.org/v07/SWH2002/index.html>
- 141 Spencer, M., Wachtel, K., & Howe, C. J. (2004). Representing Multiple Pathways of Textual
142 Flow in the Greek Manuscripts of the Letter of James Using Reduced Median Networks.
143 *Computers and the Humanities*, 38, 1–14. [https://doi.org/10.1023/B:CHUM.0000009290.](https://doi.org/10.1023/B:CHUM.0000009290.14571.59)
144 [14571.59](https://doi.org/10.1023/B:CHUM.0000009290.14571.59)
- 145 Turnbull, R. (2020). *The Textual History of Codex Sinaiticus Arabicus and Its Family*. Ridley
146 College.

DRAFT