

teiphy: A Python Package 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Published: unpublished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 in TEI XML format and can convert it to any of the following formats amenable to phylogenetic analysis: NEXUS (with support for ambiguous states and clock model calibration data blocks for MrBayes or BEAST2), Hennig86, PHYLIP (relaxed for use with RAxML), FASTA (relaxed for use with RAxML), and STEMMA (designed for Stephen C. Carlson's stemmatic software tailored for textual data). For machine learning-based analyses, teiphy can also convert a TEI XML collation to a collation matrix in NumPy, Pandas DataFrame, CSV, TSV, or Excel format.

Statement of Need

For over a decade, the Text Encoding Initiative has endeavored to provide an international standard for digitally encoding textual information for the humanities (Ide & Sperberg-McQueen, 1995). Their guidelines describe a standard format for encoding material details, textual transcriptions, and critical apparatuses (TEI Consortium, 2022). Due to its rich and well-documented set of elements for expressing a wide range of features in these settings, the Text Encoding Initiative's Extensible Markup Language format (hereafter abbreviated 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 digital transcriptions and collations 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. describe uncertainties about a source's contents as accurately as possible, allowing for degrees of uncertainty and multiple choices for disambiguations if necessary.

These principles have much bearing on the editing of critical texts, a task fundamental to the field of philology. Phylogenetic algorithms developed in the context of evolutionary biology have been popular approaches to this task, especially as philology itself has taken a digital turn over the years. 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

42 for computer-assisted textual criticism (Salemans, 2000). Over the decades preceding and
43 following these developments, biologists have continued to develop and improve phylogenetic
44 methods (Felsenstein, 2004), and textual critics have adapted these improvements and even
45 added their own innovations to make the process more suitable for their purposes (Carlson,
46 2015; Edmondson, 2019; Hyytiäinen, 2021; Spencer et al., 2002, 2004; Turnbull, 2020).

47 Phylogenetic algorithms have a natural place in textual criticism given the deep analogy between
48 textual traditions and evolutionary trees of life: a sequence alignment, which consists of taxa,
49 sites or characters, and the states of taxa at those characters, corresponds almost identically
50 to a collation, which consists of witnesses to the text, locations of textual variation (which
51 we will call “variation units” from here on), and the variant readings attested by witnesses at
52 those points.

53 The problem is that no currently available phylogenetic software accepts inputs in TEI XML
54 format. Instead, most phylogenetic programs expect inputs in NEXUS format (Maddison
55 et al., 1997). This format was conceived with versatility—including use in stemmatic and
56 non-biological applications—in mind, and this design choice has been vindicated in its general
57 applicability, but NEXUS is neither equipped nor meant to express the same kinds of details
58 that TEI XML is. Conversely, for those interested primarily in working with the collation as an
59 alignment, TEI XML is overkill. Thus, a chasm continues to separate data born in the digital
60 humanities from phylogenetic tools born in the biological sciences, and the only way to bridge
61 it is by conversion.

62 The challenge is made more daunting by the variety of tools available for phylogenetic and
63 other analyses, some of which expect inputs other than NEXUS files or NEXUS files augmented
64 in different ways. For instance, the cladistic software PAUP*, which has historically been the
65 tool of choice for text-critical applications that evaluate candidate trees using the criterion of
66 maximum parsimony, reads inputs in NEXUS format (Swofford, 2003), but the TNT software,
67 its main competitor among maximum parsimony-based programs, expects inputs in Hennig86
68 format (Farris, 1988; Goloboff & Catalano, 2016). While Hennig86 format does not allow
69 for as much flexibility in the input as NEXUS does (e.g., it does not support ambiguous
70 character states that can be disambiguated as some states but not others), TNT’s extensive
71 support for morphological state models makes it potentially more suitable for textual data, and
72 textual critics may prefer it to PAUP*. In the same regime, Stephen C. Carlson’s [STEMMA](#)
73 [software](#), initially developed for the cladistic analysis of biblical texts known to be affected by
74 contamination, makes substantial adaptations to the basic maximum-parsimony phylogenetic
75 approach to account for this problem and other constraints common in a text-critical setting
76 (Carlson, 2015); however, the input collation data must be provided in a STEMMA-specific
77 format. Likewise, among programs that use the maximum likelihood criterion instead of
78 maximum parsimony, IQ-TREE accepts inputs in NEXUS format (Minh et al., 2020), but
79 RAxML interfaces primarily with inputs in PHYLIP and FASTA format (Stamatakis, 2014).
80 Finally, for phylogenetic programs that attempt to estimate the posterior distribution of
81 candidate trees in a Bayesian fashion, MrBayes and BEAST2 both accept inputs in NEXUS
82 format (or can convert NEXUS inputs to their standard input format), but they expect taxon
83 dates (for the calibration of evolutionary clock models) to be specified in the NEXUS file in
84 different code blocks (Bouckaert et al., 2019; Ronquist et al., 2012).

85 Furthermore, end users of textual collations may be interested in non-phylogenetic analyses.
86 In this case, the desired input format is often not a NEXUS-style sequence alignment, but
87 a collation matrix with a row for each variant reading and a column for each witness (or
88 vice-versa). For Python machine-learning libraries like Scikit-learn (Pedregosa et al., 2011),
89 NIMFA (Zitnik & Zupan, 2012), and TensorFlow (Abadi et al., 2015), the standard input
90 format is a NumPy array (Harris et al., 2020), although Pandas DataFrames, which support
91 row and column labels (McKinney, 2010; The pandas development team, 2020), may also be
92 supported. (The latter format also extends the conversion pipeline to many other formats,
93 including CSV, TSV, and Excel files; Pandas DataFrames can even write their contents to
94 database tables.) To give an example, the text of the New Testament has served as a testbed

95 for multiple analyses of this type, which have generally applied clustering and biclustering
96 algorithms to collation matrices (Baldwin, 2010; Finney, 2018; McCollum, 2019; Thorpe, 2002;
97 Wilker, 2008). Given the prevalence of efforts like these, the need for a means of converting
98 TEI XML collations to NumPy collation matrices or labeled Pandas DataFrames is clear.

99 Design

100 While the conversion process is a straightforward one for most collation data, various sources
101 of ambiguity can make a one-to-one mapping of witnesses to readings impossible. One such
102 source of ambiguity is lacunae, or gaps in the text due to erasure, faded ink, or damage to the
103 page. Another is retroversions, or readings in the original language of the text reconstructed
104 through the back-translation of subsequent versions of the text in other languages. Mechanisms
105 for modeling ambiguous states resulting from situations like these exist in both TEI XML and
106 NEXUS, and in both parsimony- and likelihood-based phylogenetic methods, ambiguities about
107 the states at the leaves and even at the root of the tree can be encoded and leveraged in the
108 inference process. For these reasons, it is imperative to ensure that these types of judgments, as
109 well as other rich features from TEI XML, can be respected (and, where necessary, preserved)
110 in the conversion process.

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

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

128 Use Case

129 Due to the availability of extensive collation data for the Greek New Testament, and because
130 this project was originally developed for use with such data, we tested this library on a sample
131 collation of the book of Ephesians in thirty-eight textual witnesses (including the first-hand texts
132 of manuscripts, corrections made to manuscripts by later hands, translations to other languages,
133 and quotations from church fathers). The manuscript transcriptions used for this collation
134 were those produced by the University of Birmingham's Institute for Textual Scholarship and
135 Electronic Editing (ITSEE) for the International Greek New Testament Project (IGNTP); they
136 are freely accessible at <https://itseeweb.cal.bham.ac.uk/epistulae/XML/igntp.xml>. To achieve
137 a balance between variety and conciseness, we restricted the collation to a set of forty-two
138 variation units in Ephesians corresponding to variation units in the United Bible Societies Greek
139 New Testament (Aland et al., 2014), which highlights variation units that affect substantive
140 matters of translation.

141 In our example collation, witnesses are described in the `listWit` element under the `teiHeader`.
142 Because most New Testament witnesses are identified by numerical Gregory-Aland identifiers,

143 these witnesses are identified with @n attributes; the recommended practice is to identify such
144 elements by @xml:id attributes, but this software is designed to work with either identifying
145 attribute (preferring @xml:id if both are provided), and we have left things as they are to
146 demonstrate this feature.

147 The witness elements in the example collation also contain origDate elements that provide
148 dates or date ranges for the corresponding witnesses. Where a witness can be dated to a
149 specific year, the @when attribute is sufficient to specify this; if it can be dated within a range
150 of years, the @from and @to attributes or the @notBefore and @notAfter attributes should
151 be used; the software will work with any of these options. While such dating elements are
152 not required, our software includes them in the conversion process whenever possible. This
153 way, phylogenetic methods that employ clock models and other chronological constraints can
154 benefit from this information when it is provided.

155 Each variation unit is encoded as an app element with a unique @xml:id attribute. Within a
156 variation unit, a lem element without a @wit attribute presents the main text, and it is followed
157 by rdg elements that describe variant readings (with the first rdg duplicating the lem reading
158 and detailing its witnesses) and their attestations among the witnesses. (Situations where
159 the lem reading is not duplicated by the first rdg element, but has its own @wit attribute,
160 are also supported.) For conciseness, we use the @n attribute for each reading as a local
161 identifier; the recommended practice for readings that will be referenced elsewhere is to use
162 the @xml:id attribute, and this software will use this as the identifier if it is specified, but we
163 have only specified @xml:id attributes for rdg elements referenced in other variation units to
164 demonstrate the flexibility of the software. For witnesses with missing or ambiguous readings
165 at a given variation unit, we use the witDetail element. For ambiguous readings, we specify
166 their possible disambiguations with the @target attribute and express our degrees of certainty
167 about these disambiguations using certainty elements under the witDetail element.

168 The [TEI XML file](#) for this example is available in the example directory of the GitHub repository.
169 Full instructions for converting this file using teiphy and analyzing it with several different
170 phylogenetic packages are provided in the documentation, but here, we will walk through the
171 command-line arguments involved in installing teiphy and using it to convert our example
172 TEI XML collation (1) to a NEXUS file suitable for use with IQ-TREE, and (2) to input for
173 the STEMMA program.

174 Because teiphy is published in the [Python Package Index \(PyPI\)](#), it can be installed via the
175 command

```
176 pip install teiphy
```

177 Now we are ready to convert our TEI XML collation to a NEXUS file for IQ-TREE. Let
178 us suppose that we would like to treat reconstructions of unclear or missing text, defective
179 spellings, and orthographic variations in spelling as trivial variants for the purposes of our
180 phylogenetic analysis. We can specify this to teiphy with the -t flag for each trivial type of
181 reading as follows:

```
182 -t reconstructed -t defective -t orthographic
```

183 In addition, suppose we would like to treat placeholders for overlapping variants from larger
184 units and lacunae as missing data. We can specify this to teiphy with the -m flag for each
185 type of reading to be read as missing data as follows:

```
186 -m overlap -m lac
```

187 If, in variation units where manuscripts are corrected or have alternate readings provided by
188 other hands, our collation adds * and T suffixes to manuscript sigla to mark the work of the
189 original hand, we can tell teiphy to ignore these suffixes using the -s flag with each trivial
190 suffix:

```
191 -s "*" -s T
```

192 (Note that because the * character is reserved on the command-line, we must place it between
193 quotation marks directly after the -s flag.)

194 When corrections are made to a manuscript, they are typically sporadic, and as a result, the
195 text of corrector witnesses like 06C1 and 06C2 will tend to be too fragmentary to be useful for
196 analysis. But if we wish to assume that each corrector tacitly adopted all of the readings from
197 the previous hand that he or she did not change, then teiphy can “fill out” each corrector’s
198 text using the text of the first hand (in the case of the first corrector) or the filled-out text of
199 the previous corrector (for all subsequent correctors). Thus, 06C1 would replicate the text of
200 06* (i.e., the first hand responsible for the text of 06) where it does not introduce its own
201 readings, and 06C2 would then replicate the text of 06C1 where it does not introduce its own
202 readings. If we want to apply this transformation during the conversion process, then we can
203 specify this with the --fill-correctors flag.

204 Finally, we must specify the required arguments to teiphy, which are the input TEI XML
205 file (example/ubs_ephesians.xml) and the name of the output NEXUS file (ubs_ephesians-
206 iqtree.nexus). Note that we do not have to specify the desired output format explicitly;
207 teiphy will determine from the output filename that it should write a NEXUS file. Combining
208 the previous options and arguments, the complete command is

```
209 teiphy -t reconstructed -t defective -t orthographic -m overlap -m lac  
210 -s "*" -s T --fill-correctors  
211 example/ubs_ephesians.xml ubs_ephesians-iqtree.nexus
```

212 If we pass the resulting NEXUS file to IQ-TREE and specify appropriate settings for our textual
213 data (in this case, the Lewis Mk substitution model with ascertainment bias correction), we
214 will get an output tree like the one shown in Figure 1.

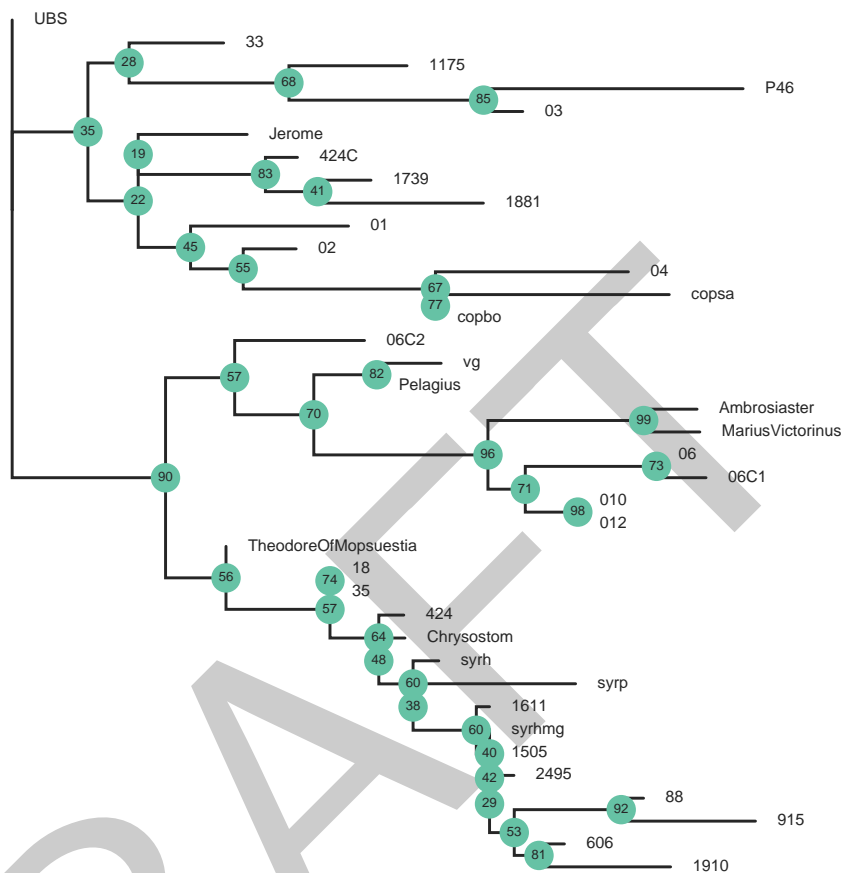


Figure 1: A phylogenetic tree inferred by IQ-TREE for the UBS Ephesians example data with support values on the branches based on 1000 bootstrap replicates. Reconstructed, defective, and orthographic sub-variants were treated as identical to their parent readings, and changes made to the text by later correctors (represented as distinct witnesses with sigla like 06C1 and 06C2) were filled in with the readings of the first hand or the previous corrector where the corrector was not active.

215 If we want to generate input files for the STEMMA program using the same options, only a few
 216 adjustments are required. First, since multiple files are written for STEMMA input (namely,
 217 a collation file with no file extension and a .chron file containing information about witness
 218 dates), we only specify the base of the filename for our output. Second, since the filename
 219 now has no extension, we must specify the desired output format to teiphy with the argument

220 `--format stemma`

221 Combining the options and arguments we used before with these changes, the complete
 222 command is

```
223 teiphy -t reconstructed -t defective -t orthographic -m overlap -m lac
224 -s "*" -s T --fill-correctors --format stemma
225 example/ubs_ephesians.xml stemma_example
```

226 If we process the output files with the PREP utility that accompanies STEMMA and then pass
 227 the resulting files to STEMMA, we will get an output tree like the one shown in Figure 2.

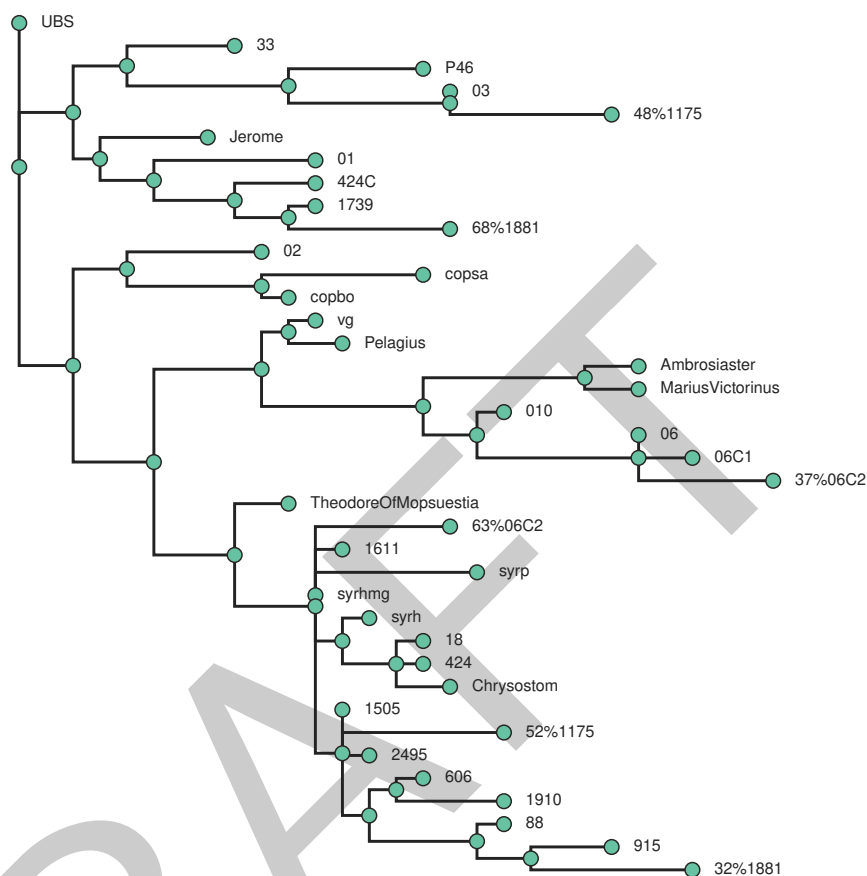


Figure 2: A phylogenetic tree inferred by STEMMA for the UBS Ephesians example data using 100 iterations of simulated annealing. Mixed witnesses are split (with proportions of their readings indicated by the percentages before their sigla) and located at different parts of the tree. Note that some witnesses (e.g., 012, 35) from the collation are excluded from this tree by STEMMA because they have the same reading sequence as another witness after their reconstructed, defective, and orthographic readings have been regularized.

228 For the small sample of variation units covered in the UBS apparatus for Ephesians, the
 229 phylogenetic results depicted in Figures 1 and 2 are impressive. The trees produced by IQ-
 230 TREE and STEMMA agree on several traditionally established groupings of manuscripts,
 231 including Family 1739 (1739, 1881, and the corrections to 424); the “Western” tradition
 232 (as preserved in the Greek-Latin diglots 06, 010, and 012, the Latin Vulgate, and the early
 233 Latin church fathers Ambrosiaster, Marius Victorinus, and Pelagius); and the later Byzantine
 234 tradition (with representative manuscripts 18 and 35 and church fathers Chrysostom and
 235 Theodore of Mopsuestia). The Harklean Syriac translation (syrh) and the witnesses to its
 236 Greek *Vorlage* (1505, 1611, 2495) are correctly placed within the Byzantine tradition, although
 237 the two programs disagree on how to describe their relationships within that tradition. While
 238 IQ-TREE does not account for mixture complicating the tradition, STEMMA identifies three
 239 witnesses suspected to exhibit Byzantine contamination: 1175, 1881, and the second corrector
 240 of 06. Both programs also identify the Codex Alexandrinus (02) as closely related to both the
 241 Sahidic and Bohairic Coptic translations of Ephesians (copsa, copbo), although they disagree
 242 on where this clade is located in the larger tradition. Despite their discrepancies regarding

243 certain subtrees, the extent of their agreements speaks to the level of genealogically significant
244 detail preserved in the TEI XML apparatus and the NEXUS and STEMMA inputs generated
245 from it.

246 Availability

247 As noted above, the software is published in [PyPI](#) and can be installed from there using `pip`.
248 The source code is available under the MIT license from the [GitHub repository](#). The automated
249 testing suite has 100% coverage, and functional tests where our example TEI XML file is
250 converted and run through RAxML, IQ-TREE, MrBayes, and STEMMA are part of `teiphy`'s
251 continuous integration (CI) pipeline.

252 Acknowledgements and Funding

253 The authors wish to thank Stephen C. Carlson for his feedback during the development of
254 `teiphy` and the *JOSS* reviewers for their thorough and insightful comments on earlier drafts of
255 this work. This work was supported by an Australian Government Research Training Program
256 (RTP) Scholarship.

257 References

- 258 Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis,
259 A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia,
260 Y., Jozefowicz, R., Kaiser, L., Kudlur, M., ... Zheng, X. (2015). *TensorFlow: Large-scale*
261 *machine learning on heterogeneous systems*. <https://www.tensorflow.org/>
- 262 Aland, B., Aland, K., Karavidopoulos, J., Martini, C. M., & Metzger, B. M. (Eds.). (2014).
263 *The Greek New Testament* (5th ed.). Deutsche Bibelgesellschaft.
- 264 Baldwin, C. S. (2010). Factor analysis: A new method for classifying New Testament Greek
265 manuscripts. *Andrews University Seminary Studies*, 48(1), 29–53.
- 266 Barbrook, A. C., Howe, C. J., Blake, N., & Robinson, P. (1998). The phylogeny of *The*
267 *Canterbury Tales*. *Nature*, 394, 839. <https://doi.org/10.1038/29667>
- 268 Bouckaert, R., Vaughan, T. G., Barido-Sottani, J., Duchêne, S., Fourment, M., Gavryushkina,
269 A., Heled, J., Jones, G., Kühnert, D., De Maio, N., Matschiner, M., Mendes, F. K.,
270 Müller, N. F., Ogilvie, H. A., du Plessis, L., Poppinga, A., Rambaut, A., Rasmussen, D.,
271 Siveroni, I., ... Drummond, A. J. (2019). BEAST 2.5: An advanced software platform
272 for Bayesian evolutionary analysis. *PLOS Computational Biology*, 15(4), 1–28. <https://doi.org/10.1371/journal.pcbi.1006650>
- 274 Carlson, S. C. (2015). *The text of Galatians and its history*. Mohr Siebeck. <https://doi.org/10.1628/978-3-16-153324-2>
- 276 Edmondson, A. C. (2019). *An analysis of the coherence-based genealogical method using*
277 *phylogenetics*. University of Birmingham. <http://etheses.bham.ac.uk/id/eprint/9150>
- 278 Farris, J. S. (1988). *Hennig86, ver. 1.5. Program and documentation*. James S. Farris.
- 279 Felsenstein, J. (2004). *Inferring phylogenies*. Sinauer Associates.
- 280 Finney, T. J. (2018). How to discover textual groups. *Digital Studies/Le Champ Numérique*,
281 8. <https://doi.org/10.16995/dscn.291>
- 282 Fischer, F. (2020). Representing the critical text. In P. Roelli (Ed.), *Handbook of stemmatology:*
283 *History, methodology, digital approaches* (pp. 405–427). De Gruyter.

- 284 Goloboff, P. A., & Catalano, S. A. (2016). TNT, version 1.5, including a full implementation
285 of phylogenetic morphometrics. *Cladistics*, 32(3), 221–238. [https://doi.org/10.1111/cla.](https://doi.org/10.1111/cla.12160)
286 [12160](https://doi.org/10.1111/cla.12160)
- 287 Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau,
288 D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van
289 Kerkwijk, M. H., Brett, M., Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., ...
290 Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, 585, 357–362. <https://doi.org/10.1038/s41586-020-2649-2>
291
- 292 Hyttiäinen, P. (2021). The changing text of Acts: A phylogenetic approach. *TC: A Journal of*
293 *Biblical Textual Criticism*, 26, 1–28.
- 294 Ide, N. M., & Sperberg-McQueen, C. M. (1995). The TEI: History, goals and future. *Computers*
295 *and the Humanities*, 29(1), 5–15. https://doi.org/10.1007/978-94-011-0325-1_2
- 296 Maddison, D. R., Swofford, D. L., & Maddison, W. P. (1997). NEXUS: An extensible file
297 format for systematic information. *Systematic Biology*, 46(4), 590–621. [https://doi.org/](https://doi.org/10.1093/sysbio/46.4.590)
298 [10.1093/sysbio/46.4.590](https://doi.org/10.1093/sysbio/46.4.590)
- 299 McCollum, J. (2019). Biclustering readings and manuscripts via non-negative matrix factoriza-
300 tion, with application to the text of Jude. *Andrews University Seminary Studies*, 57(1),
301 61–89.
- 302 McKinney, W. (2010). Data structures for statistical computing in Python. In S. van der Walt
303 & J. Millman (Eds.), *Proceedings of the 9th Python in Science Conference* (pp. 56–61).
304 <https://doi.org/10.25080/Majora-92bf1922-00a>
- 305 Minh, B. Q., Schmidt, H. A., Chernomor, O., Schrempf, D., Woodhams, M. D., von Haeseler,
306 A., & Lanfear, R. (2020). IQ-TREE 2: New models and efficient methods for phylogenetic
307 inference in the genomic era. *Molecular Biology and Evolution*, 37(5), 1530–1534. <https://doi.org/10.1093/molbev/msaa015>
308
- 309 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
310 Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D.,
311 Brucher, M., Perrot, M., & Duchesnay, É. (2011). Scikit-learn: Machine learning in Python.
312 *Journal of Machine Learning Research*, 12(85), 2825–2830.
- 313 Robinson, P., & O'Hara, R. J. (1992). Report on the Textual Criticism Challenge 1991. *Bryn*
314 *Mawr Classical Review*, 3(4), 331–337.
- 315 Ronquist, F., Teslenko, M., van der Mark, P., Ayres, D. L., Darling, A., Höhna, S., Larget, B.,
316 Liu, L., Suchard, M. A., & Huelsenbeck, J. P. (2012). MRBAYES 3.2: Efficient Bayesian
317 phylogenetic inference and model selection across a large model space. *Systematic Biology*,
318 61(3), 539–542. <https://doi.org/10.1093/sysbio/sys029>
- 319 Salemans, B. J. P. (2000). *Building stemmas with the computer in a cladistic, neo-Lachmannian,*
320 *way: The case of fourteen text versions of Lanceloet van Denemerken*. Katholieke Univer-
321 siteit Nijmegen. <https://hdl.handle.net/2066/147058>
- 322 Spencer, M., Wachtel, K., & Howe, C. J. (2002). The Greek Vorlage of the Syra Harclensis:
323 A comparative study on method in exploring textual genealogy. *TC: A Journal of Biblical*
324 *Textual Criticism*, 7. <http://jbtc.org/v07/SWH2002/index.html>
- 325 Spencer, M., Wachtel, K., & Howe, C. J. (2004). Representing multiple pathways of textual flow
326 in the Greek manuscripts of the Letter of James using reduced median networks. *Computers*
327 *and the Humanities*, 38, 1–14. <https://doi.org/10.1023/B:CHUM.0000009290.14571.59>
- 328 Stamatakis, A. (2014). RAxML version 8: A tool for phylogenetic analysis and post-
329 analysis of large phylogenies. *Bioinformatics*, 30(9), 1312–1313. [https://doi.org/10.](https://doi.org/10.1093/bioinformatics/btu033)
330 [1093/bioinformatics/btu033](https://doi.org/10.1093/bioinformatics/btu033)

- 331 Swofford, D. L. (2003). *PAUP*: Phylogenetic analysis using parsimony (*and other methods).*
332 *Version 4.* Sinauer Associates.
- 333 TEI Consortium. (2022). *TEI P5: Guidelines for electronic text encoding and interchange:*
334 *Critical apparatus [v.4.4.0].* <https://www.tei-c.org/release/doc/tei-p5-doc/en/html/TC.html>.
335
- 336 The pandas development team. (2020). *Pandas-dev/pandas: pandas.* Zenodo. <https://doi.org/10.5281/zenodo.3509134>
337
- 338 Thorpe, J. C. (2002). Multivariate statistical analysis for manuscript classification. *TC: A*
339 *Journal of Biblical Textual Criticism*, 7. <http://jbtc.org/v07/Thorpe2002.html>
- 340 Turnbull, R. (2020). *The textual history of Codex Sinaiticus Arabicus and its family.* Ridley
341 College.
- 342 Willker, W. (2008). *Principal component analysis of manuscripts of the Gospel of John.*
343 <http://www.willker.de/wie/TCG/PCA/index.html>
- 344 Zitnik, M., & Zupan, B. (2012). NIMFA: A Python library for nonnegative matrix factorization.
345 *Journal of Machine Learning Research*, 13, 849–853.

DRAFT