

White Paper

Report ID: 112322

Application Number: HD-51907-14

Project Director: Caroline T. Schroeder (cschroeder@pacific.edu)

Institution: University of the Pacific

Reporting Period: 5/1/2104-12/31/2016

Report Due: 12/31/2016

Date Submitted: 12/30/2016

Coptic SCRIPTORIUM: A Corpus, Tools, and Methods for Corpus Linguistics and Computational Historical Research in Ancient Egypt

White Paper

National Endowment for the Humanities Office of Digital Humanities

Grant HD-51907-14

Director: Caroline T. Schroeder, University of the Pacific

Co-Director: Amir Zeldes, Georgetown University

Project website: copticscriptorium.org

December 30, 2016

Table of Contents

[1. Project Description](#)

[2. Accomplishments and Grant Products](#)

[2.1 Digitized Coptic Corpus in Multiple Formats and Visualizations](#)

[2.2 Digital and Computational Tools](#)

[2.3 ANNIS Database](#)

[2.4 Collaborative Platform](#)

[2.5 Documentation](#)

[2.6 Web application to provide easier user access for reading and citing visualizations of textual data](#)

[2.7 Symposium and workshop \("Digital Coptic 2." March 2015\)](#)

[2.8 Public tutorial and workshop](#)

[2.9 Articles and conference papers to distribute the results of our work](#)

[Peer-reviewed Articles](#)

[Published Conference Proceedings](#)

[Other Conference Presentations and Lectures](#)

[3. Project Activities](#)

[4. Audiences and Impact](#)

[5. Evaluation](#)

[6. Continuation of the Project](#)

[Appendices](#)

[Appendix A: List of Participants](#)

[Appendix B: Website Screenshots](#)

[Appendix C: Data Visualizations](#)

[Appendix D: Digital Coptic 2 Program](#)

[Appendix E: Guidelines for Citation, Transcription, Part of Speech Tagging, Lemmatization](#)

1. Project Description

Coptic, having evolved from the language of the hieroglyphs of the pharaonic era, represents the last phase of the Egyptian language and is pivotal for a wide range of disciplines, such as linguistics, biblical studies, the history of Christianity, Egyptology, and ancient history. Coptic SCRIPTORIUM has developed the first open-source technologies for computational and digital research across the disciplines as applied to Egyptian texts. We created a digitized corpus of Coptic texts available in multiple formats and visualizations (including TEI XML), tools to analyze and process the language (e.g., the first Coptic part-of-speech tagger), a database with search and visualization capabilities, and a collaborative platform for scholars to contribute texts and annotations and to conduct research. The technologies and corpus function as a collaborative environment for digital research by any scholars working in Coptic.

During the funded period, the project also received support from an NEH Preservation and Access Division, Humanities Collections and Reference Resources Foundations Grant (PW-51672-14). This White Paper contains a holistic account of project activities, concentrating on work pertaining to the ODH grant while also including activities that intersected with our other grant. This White Paper may contain similar, or in some places identical, phrasing as in the Final Performance Report and White Paper for PW-51672-14.¹

2. Accomplishments and Grant Products

The project achieved the primary objections we set out to accomplish in the original proposal:

1. a digitized corpus of Coptic texts available in multiple formats and visualizations
2. digital and computational tools to analyze, process, and visualize the language (e.g., the first Coptic part-of-speech tagger, converters between digital formats and fonts, etc.)
3. a database created from the texts and tools using the ANNIS search and visualization infrastructure
4. a collaborative platform for scholars to contribute texts and annotations as well as conduct research using the corpus, tools, and database
5. documentation of project tools and technology

¹ Schroeder, Caroline T., Zeldes, Amir, and Elizabeth Platte. *NEH White Paper Report for Coptic SCRIPTORIUM: Digitizing a Corpus for Interdisciplinary Research in Ancient Egyptian*. National Endowment for the Humanities, 2016. Search for grant, White Paper, and grant products at: <https://securegrants.neh.gov/publicquery/main.aspx>. Also available on our website at <http://copticSCRIPTORIUM.org/download/Coptic-SCRIPTORIUM-white-paper-2016-NEH-PW-51672-14.pdf>.

In addition to these outcomes, the project also achieved the following accomplishments, which were partially co-funded by an NEH Preservation and Access, Humanities Collection and Reference Resources Foundation grant (PW-51672-14):

6. Web application to provide easier user access for reading and citing visualizations of textual data
7. Symposium and workshop on digital Coptic Studies (“Digital Coptic 2,” March 2015)

Finally, the project also succeed in accomplishing:

8. Public tutorial and workshop on Coptic SCRIPTORIUM at the International Association of Coptic Studies
9. Articles and conference papers to distribute the results of our work

2.1 Digitized Coptic Corpus in Multiple Formats and Visualizations

Coptic SCRIPTORIUM uses the tools described in 2.2 to annotate digitized Coptic texts. Due to our interdisciplinary approach, we release our data in three formats.

- XML files using the EpiDoc subset (<https://sourceforge.net/p/epidoc/wiki/Home/>) of the Text Encoding Initiative (TEI) XML standards (<http://www.tei-c.org/index.xml>). EpiDoc is the standard in the field for producing digital editions of epigraphy, ancient manuscripts, papyrology, and other ancient text-bearing objects. Our EpiDoc TEI XML files do not contain the full set of annotations available in our corpora. The annotations in these files include information about text and manuscript structure, core philological and linguistic annotations (such as part of speech and loan words), and most metadata.
- XML files in standoff annotation using the PAULA XML format (<https://hal.inria.fr/hal-00783716>). These files contain the complete dataset of text and annotations for all corpora.
- Relational database files for use in the ANNIS search and visualization infrastructure (<http://corpus-tools.org/annis/>). ANNIS is the web-based database Coptic SCRIPTORIUM uses for search and visualization of the corpora. ANNIS can also be installed on a user’s desktop; the researcher can load our relANNIS database files into this local installation.

All the files are released under a Creative Commons Attribution (<https://creativecommons.org/licenses/by/4.0/>) license (either CC-BY 3.0 or 4.0 depending on date of corpus release) except the corpora derived from the Sahidica project, which operate under Sahidica’s original more limited license for academic use only.

All files can be downloaded from our Github site’s corpus repository at github.com/CopticScriptorium/corpora. We provide many links to this GitHub repository throughout our site. The version control system inherent in Git allows researchers to access previous versions of the digital files even as we update the corpora.

We have two web-based services to read and query the texts and annotations. On a Georgetown University server, an instance of the ANNIS database provides searchable access to our corpus and annotations (see [section 2.3](#) below). We also developed in collaboration with the concurrent PW-51672-14 grant a web application to enable easier access to and citation of texts using the CTS (Canonical Text Services) citation model, available at <http://data.copticscriptorium.org>.² (See also [section 2.6](#) below.)

Currently in addition to the search capacity in ANNIS, the annotated text data is automatically serialized into a reader friendly HTML format and styled using several CSS stylesheets to produce multiple visualizations in HTML. The visualizations are generated dynamically and cached, meaning that updates to the corpus can easily be made browsable, but access to visualizations is instantaneous for readers. Potential visualizations are expandable based on the data model. Each visualization is dependent on various combinations of annotations in our data model (which is described in the following two sections). Current visualizations include:

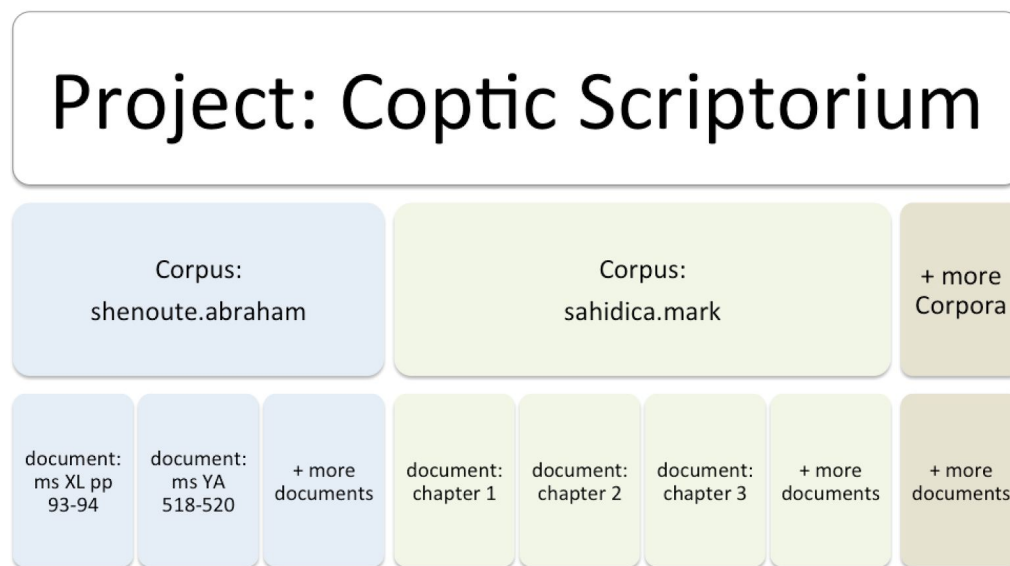
- Normalized view: the normalized Coptic text segmented into bound groups, with an English translation (when available) appearing as a pop-up when a cursor hovers over the text. Optimal for people wishing to read Coptic. Under the aegis of the NEH-DFG grant for the KELLIA project (HG-229371, see <http://copticscriptorium.org/kellia>), the normalized words now appear as hyperlinks to entries in a new online Coptic Dictionary (at <https://corpling.uis.georgetown.edu/coptic-dictionary>). Under the KELLIA grant, work is ongoing to enable proper mapping of words in our segmentation system to words segmented according to our German partners' system. There are currently occasional words where the links result in an empty search in the Dictionary.
- Diplomatic view: the diplomatic transcription of a manuscript page, which resembles the appearance of the manuscript page.
- Analytic view: an aligned visualization of the normalized Coptic text, part of speech tags, and the English translation. The manually curated Biblical corpora (currently the Gospel of Mark and 1 Corinthians) include an alignment with the Greek Bible, provided by the Society of Biblical Literature and Logos Software.
- Chapter view: Similar to the normalized view but divided into numbered chapters and verses for corpora such as the Bible which have existing canonical versification. As with the Normalized view, words link to the online Coptic dictionary.
- Histogram visualization for aggregate queries (frequency analysis for combinations of features extracted from a query).

Please see Appendix C for examples of these visualizations.

All corpora are developed and released using version control. Since this is the first project to develop natural language processing tools for the Egyptian language family, the tools and data

² Almas, Bridget, and Caroline T. Schroeder. "Applying the Canonical Text Services Model to the Coptic SCRIPTORIUM." *Data Science Journal* 15.0 (2016): p.13. doi:10.5334/dsj-2016-013 (<http://datascience.codata.org/articles/10.5334/dsj-2016-013/>)

model have evolved over the course of the project. (For more on the multilayer data model, see the White Paper for the NEH Preservation and Access grant PW-51672-14.) Our corpus architecture consists of one or more documents within each corpus; each corpus is a logical collection of documents based on existing research standards (e.g., the shenoute.abraham corpus consists of all documents pertaining to the text by Shenoute known as *Abraham Our Father*, and the sahidica.mark corpus consists of all documents pertaining to the Sahidica online version of the Gospel of Mark). The Coptic SCRIPTORIUM project then consists of multiple corpora, with multiple documents within each corpora, as demonstrated in the example figure below:



Coptic Scriptorium Corpus Architecture

This structure is motivated in part by the functionality of the ANNIS search and visualization tool we chose to use for a database and in part by the process of digitizing and annotating the corpora in parts. We also have been working in a public, open access environment, meaning that we have released our corpora as documents have become available. We release a new version of a corpus when one or more new documents in that corpus have been digitized and annotated, or when one or more documents in that corpus have been edited due to updates in the data model, updates to the tools, or corrections to the annotations. Developing a version numbering system for a complex corpus architecture was one of the challenges of the project. After some experimentation, we now follow a system in which a new or newly edited document receives a version number which is the same as the new version number of the corpus, which also is the same as number given to the overall the project release. The overall project release consists of an archive file of all corpora files (in the three data formats described above) released on our GitHub repository at <https://github.com/CopticScriptorium/corpora/releases>.

2.2 Digital and Computational Tools

For machine processing and annotation of digital Coptic text, the project produced the following tools and technologies. The tools were developed for the Sahidic dialect of the Coptic language. Development took place primarily under the aegis of this grant (HD-51907), although the editors and annotators of corpora, who contributed data that increased the accuracy of the tools, were funded under the concurrent grant PW-51672-14. Editors used them in conjunction with manual annotation to produce the annotated text corpora (see [section 2.1](#)). All tools are available under open licenses. Documentation for these tools can be found in their respective GitHub repositories, linked below. Detailed guidelines for transcription (including tokenization), part of speech tagging, and lemmatization appear in Appendix E.

- Font and character converters: Convert text of documents transcribed in legacy ASCII fonts into the Unicode Coptic character set. Project participants produced a converter for several legacy fonts (`recode_coptic`) and one for a more complex legacy font created originally by Dirk Van Damme and Gregor Wurst (`CopticVDWtoUTFConverter`). Collaborator So Miyagawa produced another (`copticizer`) for text in the St. Shenouda Society CDs. All converters are on the Converters repository on GitHub (<https://github.com/CopticScriptorium/converters>).
- Tokenizer: Segments Coptic bound groups into words and morphs. This tool requires text input segmented according to the principles of boundedness articulated in Bentley Layton's *Coptic Grammar*.³ Available on the Tokenizers repository on GitHub (<https://github.com/CopticScriptorium/tokenizers>).
- Normalizer: Normalizes orthography and spelling of Coptic text. Available on the Normalizer repository on GitHub (<https://github.com/CopticScriptorium/normalizer>).
- Part of speech tagger: A probabilistic tagger built from the independent tool TreeTagger trained on a set of Coptic training data. The most recent version includes lemmatization data as well as part of speech annotation, and the lemma information includes Greek Loanword lemmas from the Dictionary and Database of Greek Loanwords in Coptic project (<http://research.uni-leipzig.de/ddglc/>). The original version was created with the assistance of a lexicon provided by the Corpus dei Manoscritti Copti Letterari project (<http://www.cmcl.it/>). Available on the Part of Speech Tagger repository on GitHub (<https://github.com/CopticScriptorium/tagger-part-of-speech>).
- Lemmatizer: A lexical tagger that annotates each word with its lemma, or dictionary headword. Originally developed as a stand alone tool and now incorporated into the part of speech tagger (above).
- Language of origin tagger: A lexical tagger that annotates each word or morph for its language of origin. (So if a text corpus has been segmented or tokenized to account for compounds, this tagger will annotate the loanword within the compound.) Native Egyptian vocabulary remain un-annotated.

³ Layton, Bentley. *A Coptic Grammar*. 3rd Edition, Rev. Wiesbaden: Harrassowitz, 2011. Print. *Porta Linguarum Orientalium Neue Serie* 20.

- Format converters:
 - SGML input/output plugin: Converts digitized text annotated with SGML or XML tags into a multilayer spreadsheet format in Microsoft Excel. Available on GitHub (<https://github.com/CopticScriptorium/converters>).
 - Annotation validation plugin: Validates multilayer annotations to ensure they conform to the data model described above. Available on GitHub (<https://github.com/CopticScriptorium/XLAddIns>).
 - EpiDoc TEI conversion tool: Converts select annotations of a multilayer text and annotations spreadsheet document into a single EpiDoc TEI-XML file. Currently operates only for Windows. Available on GitHub (<https://github.com/CopticScriptorium/XLAddIns>).

Over the course of the grant period, we updated the tools with periodic enhancements and bug fixes. In addition, we recursively updated the tools after publishing a large amount of new textual data. Since editors had manually edited the annotations, correcting any errors or omissions of the tools, using published data to update the tools resulted in greater accuracy.

Under the aegis of a bilateral NEH-DFG grant (HG-229371) of the KELLIA collaboration with German partners, the project also updated the tools with a Greek lemma list created by the [Database and Dictionary of Greek Loanwords in Coptic](#) project. In addition, funding from this grant supported development of a natural language processing pipeline, which runs any or all of the tools in a web application with an API.⁴ The pipeline can be used at <https://corpling.uis.georgetown.edu/coptic-nlp/>, and the code is available on GitHub (<https://github.com/CopticScriptorium/coptic-nlp>).

2.3 ANNIS Database

ANNIS is an open source, browser based search and visualization platform for linguistic corpora, which we have re-purposed and adapted in places in order to grant intuitive and flexible access to the annotations underlying our data model and visualizations. During the current project period, several improvements were made to the infrastructure to support our project, in particular:

- Embeddable visualization links to grant easy access to views of whole documents from external applications, including our web application for close reading purposes (rather than quantitative searches) at data.copticscriptorium.org
- Support for HTML templates to customize the appearance of specific annotation layers within generic visualizations

⁴ Amir Zeldes and Caroline T. Schroeder. “An NLP Pipeline for Coptic.” In: *Proceedings of the 10th ACL SIGHUM Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities (LaTeCH2016)*. Berlin, 146–155. doi:10.18653/v1/W16-2119.

- Hyperlinking in both document and grid views. This feature is now used to link search results and archived ANNIS corpora in the repository to the live Coptic Dictionary Online at <https://corpling.uis.georgetown.edu/coptic-dictionary/>

The multilayer data model and the metadata model used in ANNIS for our corpora were partially developed under Preservation and Access Grant PW-51672-14; see the White Paper section 2.2 for more details.⁵

2.4 Collaborative Platform

Coptic SCRIPTORIUM's GitHub repository and website have developed into a collaborative platform where project editors contribute transcriptions and annotations to the project. Annotators/editors funded by Grant PW-51672-14 forked the corpora, and then used the tools developed under this ODH Startup Grant to annotate the corpora. In addition participants and other collaborators have contributed to the development of the tools through the GitHub site.

Project co-PIs and editors have identified the need for a text transcription editor (with version controlled commits to a server, such as a university server or a GitHub repository) and a multi-layer editing tool (also with version controlled commits to a server) to enable annotations and editing in the multilayer model. Inline editing and annotation with our complex data model is not feasible. We are pursuing the development of these two tools under the NEH-DFG grant for the KELLIA project (HG-229371). (See also [section 6](#) Continuation.)

2.5 Documentation

Documentation, including licensing information, for each tool appears in the GitHub repository with the tool, either in a separate file for Guidelines or in the ReadMe document for the repository. We also have created a wiki for project data models, instructions, and workflows at wiki.copticscriptorium.org, and a blog announcing news, corpus and tool releases, and other information at blog.copticscriptorium.org. We have a YouTube channel with video tutorials (<https://www.youtube.com/playlist?list=PLDN6-yDloRwf5HNCwZOAbbSfih1iEqOJT>). Our project website includes a tutorial with sample queries for the ANNIS database (<http://copticscriptorium.org/ANNIS-tips.html>), and a page with instructions for how to cite our work (<http://copticscriptorium.org/citation-guidelines.html>); also see Appendix E for our Citation Guidelines. In addition, all NEH project reports appear on a reports page of the website (<http://copticscriptorium.org/reports>). All of these resources can be accessed easily from links on our main Documentation page. The documentation was produced under this grant as well as the concurrent NEH-DFG grant for the KELLIA project (HG-229371), and the NEH Preservation and Access, HCRR Foundation grant (PW-51672-14).

⁵ Schroeder, Zeldes, and Platte, *op. cit.*

2.6 Web application to provide easier user access for reading and citing visualizations of textual data

We developed a web application that provides the most recent version of our data in all the formats currently available using the CTS (Canonical Text Services) URN system. Researchers can enter the the URN for a document or set of documents and retrieve the latest version(s). When the URN resolves, the web service returns with links to all all visualizations for the relevant document(s), link(s) to the ANNIS search tool, and links to data downloads in all available formats (PAULA, TEI, relANNIS) on GitHub. The web service uses the ANNIS API to present the data in manner more accessible to researchers who wish to read or browse documents rather than search across corpora. Code is deployed at <http://data.copticscriptorium.org> and released open source under Apache 2.0 and Creative Commons Attribution 4.0 (CC-BY 4.0) licenses on GitHub (<https://github.com/CopticScriptorium/cts>). Funding from this grant as well as the NEH Preservation and Access grant PW-51672-14 supported its development. See Appendix B for screenshots.

2.7 Symposium and workshop (“Digital Coptic 2,” March 2015)

The project hosted a two-day symposium and workshop at Georgetown University in March 2015. On day one, scholars in Coptic Studies and Digital Humanities presented papers on ongoing digital projects or research using digital and computational methods. During discussion sessions, participants focused on how to collaborate across projects and technologies. On day two, Coptic SCRIPTORIUM participants led a tutorial on how to use our project resources; we also solicited feedback from the participants about the future direction of the project and suggestions generally. Many of the presentations plus part of the tutorial are available on our [YouTube channel](#). This event was co-funded by an NEH Preservation and Access, Humanities Collection and Reference Resources Foundation grant (PW-51672-14). The program and schedule appear in Appendix D.

2.8 Public tutorial and workshop

At the 2016 Congress of the International Association of Coptic Studies, held at Claremont Graduate University in July 2016, we held a tutorial and workshop about the project. Project Participants Amir Zeldes, Caroline T. Schroeder, and Rebecca Krawiec directed it. We introduced participants to our natural language processing tools, to searching and querying in the ANNIS database, and to reading and citing documents. We also received feedback from participants regarding the project.

2.9 Articles and conference papers to distribute the results of our work

The project produced the following articles and conference papers during the grant period. Some of these articles and papers document research conducted under the aegis of the NEH Preservation and Access grant PW-51672-14. For completeness and accuracy, we include all the papers and articles produced during the grant period, even if research was supported in part by the other grant.

Peer-reviewed Articles

Almas, Bridget, and Caroline T. Schroeder. "Applying the Canonical Text Services Model to the Coptic SCRIPTORIUM." *Data Science Journal* 15 (2016): p.13. doi:10.5334/dsj-2016-013 ([link](#))

Schroeder, Caroline T., and Amir Zeldes. "Raiders of the Lost Corpus." *Digital Humanities Quarterly* 10.2 (2016). ([link](#))

Zeldes, Amir, and Caroline T. Schroeder. "Computational Methods for Coptic: Developing and Using Part-of-Speech Tagging for Digital Scholarship in the Humanities." *Digital Scholarship in the Humanities* 30.suppl 1 (2015): i164–i176. doi:10.1093/llc/fqv043. ([link](#))

Published Conference Proceedings

Zeldes, Amir. "Duplicitous Diabolos: Parallel witness encoding in quantitative studies of Coptic manuscripts." Presented at Symposium on Cultural Heritage Markup, Washington, DC, August 10, 2015. In *Proceedings of the Symposium on Cultural Heritage Markup*. Balisage Series on Markup Technologies, vol. 16 (2015). doi:10.4242/BalisageVol16.Zeldes01. ([link](#))

Schroeder, Caroline T. "Shenoute in Code: Digitizing Coptic Cultural Heritage for Collaborative Online Research and Study." *Coptica* 14 (2015): 21-36. ([link](#))

Other Conference Presentations and Lectures

Schroeder, Caroline T. "The Future of Biblical Scholarship in a Digital Age." Plenary Presentation for the Catholic Biblical Association. Santa Clara. August 2016.

-----, "Coptic SCRIPTORIUM: A Digital Platform for Research in Coptic Language and Literature." Congress of the International Association of Coptic Studies. Claremont, Ca. July 2016.

Krawiec, Rebecca S. "Charting Rhetorical Choices in Shenoute: *Abraham our Father* and *I See Your Eagerness* as case-studies." Congress of the International Association of Coptic Studies. Claremont, Ca. July 2016.

- Luckritz Marquis, Christine. "Reimagining the *Apophthegmata Patrum* in a Digital Culture." Congress of the International Association of Coptic Studies. Claremont, Ca. July 2016.
- Zeldes, Amir. "A Quantitative Approach to Syntactic Alternations in Sahidic." Congress of the International Association of Coptic Studies. Claremont, Ca. July 2016.
- Schroeder, Caroline T. "Preserving Coptic Cultural Heritage for the Digital Future," Seventeenth St. Shenouda-UCLA Conference of Coptic Studies. Los Angeles. July 17-18, 2015.
- Platte, Beth. "Coptic SCRIPTORIUM: Data from the Desert," Linking the Big Ancient Mediterranean. University of Iowa. June 6-8, 2016.
- Zeldes, Amir. "Tagging the Desert Fathers: Part of Speech Analysis in Sahidic Coptic Corpora." 43rd Annual North American Conference on Afroasiatic Linguistics (NACAL2015), 13-15 February 2015. Washington, DC.
- Schroeder, Caroline T. "Tag, You're It: Creating a Richly Annotated Coptic Digital Library," Society of Biblical Literature Annual Meeting. San Diego. November 2014
- , "DH Technologies for the Study of Coptic Language and Literature," Brown University Library, Digital Lab, September 30, 2014.
- Schroeder, Caroline T. and Amir Zeldes. "Digitizing the Dead and Dismembered: DH Technologies for the Study of Coptic Texts." DH2014. Lausanne, Switzerland. July 2014.
- , "Tagging Shenoute." North American Patristics Society Annual Meeting. Chicago. May 2014.
- , "Digital Coptic: Building an Online Environment for the Study of Coptic Literature." Center for Tebtunis Papyri, University of California, Berkeley. May 2014.

3. Project Activities

The primary project activities were:

- Development of tools described in [section 2.2](#)
- Application of tools to develop the annotated corpora (partially conducted by editors on the concurrent grant PW-51672-14), see [section 2.1](#) and [section 2.3](#). Editors or senior editors transcribed and/or annotated and/or translated the text. A senior editor reviewed and edited each document prior to publication to ensure editorial review.
- Software development on the document web application described in section by contractors (Archimedes and DaveBSoft).
- Digital Coptic 2 Workshop
- Writing Documentation

- Project management: coordinating editors, software development, publicity for project, budget, and communication between project participants
- Conducting research and writing articles and conference papers (see [section 2.9](#))
- Project meetings in San Francisco, CA (May 2014); Washington DC (March, 2015); Stockton, CA (May 2015); Washington DC (December 2015); Claremont CA (July 2016); Washington DC (December 2016). Some of the travel for these meetings was funded by the PW-51672-14 and HG-229371 grants.
- Advisory Board meetings conducted over Skype and email 2-3 times per year. Additionally, project co-directors consulted with Board members individually regularly throughout the grant period.

Project results were publicized at conferences and in articles (see [section 2.9](#)), on our project website and blog, and on social media.

4. Audiences and Impact

The audience consists primarily of academics (faculty, researchers, and graduate students) in Linguistics, Egyptology, and Religious Studies. These groups are the most prominent audiences at presentations and conferences and have provided the most anecdotal feedback to project directors. Additionally, undergraduates studying Coptic have used the site, according to anecdotal reports from faculty at other institutions. Social media interactions also indicate that non-academics with interest in Coptic, as well as members of the Coptic Orthodox laity in the American and European diaspora, visit the project site and follow its progress.

Our corpora have been forked from GitHub and/or used in research by several scholars. Four researchers or projects have forked our corpora, including the Classical Language Toolkit (an aggregator of natural language processing tools and corpora for ancient languages). Paul Dilley (University of Iowa) and So Miyagawa (graduate student, U. of Göttingen) used corpora and tools for their papers at our March 2015 symposium and at the International Association of Coptic Studies. The [Text Alignment Network](#) uses our edition of the Coptic New Testament in its online, aligned six-language New Testament. Rebecca S. Krawiec has used our corpora in her forthcoming article “Reading Abraham in the White Monastery: Fathers, Sources, and History,” to appear in [REDACTED]

According to Google Analytics, our main web domain copticscriptorium.org has been visited by over 12,000 unique users since it was launched in 2014, with just over 5000 repeat visitors. Most site views come from the United States; other significant audiences are in Germany, the UK, Brazil, Egypt, Japan, Canada, and Russia. The project blog, launched just over a year ago, has had over 2000 unique visitors. By comparison, approximately 200 scholars attended the 2016 Congress for the International Association of Coptic Studies.

5. Evaluation

The project was continually evaluated internally and periodically externally. In email correspondence and at project meetings project participants evaluated our progress and outputs. Advisory board meetings were held multiple times a year, over email and via Skype, since board members live in 3 different time zones. Feedback from our board meetings was incorporated into the project. We also consulted with board members individually about specific issues regarding linguistics, data curation, sharing data with other projects (such as papyri.info and the DDGLC), and translation. At two workshops (in Washington DC in 2015 and in Claremont, CA in 2016) we solicited comments, suggestion, and feedback. We also welcome feedback via social media, email to our project, and GitHub issues and pull requests.

The three primary concerns expressed by outside users were 1) a need for more digitized, annotated corpora in order to conduct more research; 2) greater accuracy of our tools; 3) difficulty of use of tools for the scholarly community without deep programming skills. Project participants agree with this assessment. Since this grant funded the creation of the tools and a small start-up corpus, however, we believe we have achieved the grant's objectives, and that addressing concerns #1 and #2 will take significantly more resources (both in time and funding). We have begun to address #3 by providing more documentation, YouTube video tutorials, and developing a web application for the natural language processing toolset (an online NLP pipeline discussed at the end of section 2.1 above). The NLP pipeline makes the toolset more accessible via a web interface and an API; it nonetheless still requires some technical skill to apply to one's research. This latter issue is a tension in many DH projects — the expectation that digital tools will make basic research processes faster, more accessible vs digital tools will enable more complex research (and thus can be complex to implement). We are currently looking at multiple models for applying the tools to develop more digitized corpora. Whereas previously, a single editor would both digitize and annotate a single text document, we are now looking at workflows that would allow some contributors to digitize text and others to use the tools to annotate. We hope this will allow us to move more quickly toward the goal of an expanded corpora. We also are pursuing funding for expanding the corpora and refining the tools.

6. Continuation of the Project

Now that we have built the core tools and a pilot corpus for digital Coptic studies, we are well-positioned to expand the corpus and refine the toolset. In particular, we intend to seek support for and continue preliminary work on the Coptic Treebank, which will enable high quality automatic parsing. This development will be of great value in itself, since it will allow us to add syntactic analysis to all of our resources, which in turn gives us access to questions like 'which verbs is a particular noun the subject or object of?' (e.g. finding all actions carried out by 'the

Devil', or all modifiers of 'blasphemy', etc.), and allows the linguistic examination of Coptic syntax across authors, genres, time and in the future, potentially also dialects.

Beyond querying syntactic information per se, we plan to use parsed data as a gateway technology to the development of Named Entity Recognition (NER) and Entity Linking for Coptic. This development will allow us to find occurrences of the same people and places mentioned in different ways, find entities across documents, and link our resources to other projects handling entities that overlap with ours.

Both the project PIs (Schroeder and Zeldes) will continue to work on the project as part of their ongoing research. Other participants, such as Rebecca Krawiec (Canisius College), and Christine Luckritz Marquis (Union Presbyterian Seminary) will as well.

Our tool development led to the aforementioned successful joint NEH-DFG grant proposal so that German digital Coptic scholars could implement our tools for annotating their texts. Our Digital Coptic workshop in our pre-NEH funded phase and our second symposium and workshop in March 2015 brought together scholars from other institutions and projects to discuss collaborations, standards, and shared challenges. The conversations initiated at these meetings led to the KELLIA project, which received funding for Göttingen, Münster, the Berlin-Brandenburg Academy, and other German partners to collaborate with us and scholars at the University of Minnesota and University of Iowa.

Appendices

Appendix A: List of Participants

Caroline T. Schroeder, the University of the Pacific

Amir Zeldes, Georgetown University

Elizabeth Platte, Reed College, Digital Humanities Specialist and Project Manager (2015-); editor and encoder/annotator (2013-)

Rebecca S. Krawiec, Canisius College, senior editor and encoder/annotator, translator (2013-)

Christine Luckritz Marquis, senior editor and encoder/annotator, translator (2014-)

So Miyagawa, University of Göttingen, editor and encoder/annotator, translator (2014-)

Elizabeth Davidson, Southern School of Energy and Sustainability, editor and encoder/annotator (2015-)

Dana Robinson, Creighton University, editor and encoder/annotator, translator (2016-)

Shuo Zhang, architecture and infrastructure (2015-)

Emma Manning, architecture and infrastructure (2016-)

Dave Briccetti, programmer and consultant (2015-)

Anthony Alcock, University of Kassel, translator (2015)

Eliese-Sophia Lincke, Humboldt University (2014-)

Bridget Almas, Perseus Digital Library, consultant for SCRIPTORIUM (2014-)

David Sriboonreuang, University of the Pacific student, intern and project manager (2015)

Lauren McDermott, the University of the Pacific student; TEI encoder and HTML programmer (2013-14)

Janet Timbie, the Catholic University of America, editor and annotator (2013)

Luke Hollis, Archimedes Digital, consultant and programmer for canonical referencing system (2014-2015)

Yanrui Liu, M.A., University of the Pacific, repository and website management (2014-2015)

Edwin Ko, Georgetown University, annotation interface development (2014)

Alex Dickerson, the University of the Pacific, student; TEI encoder and programmer (2013)

Advisory Board:

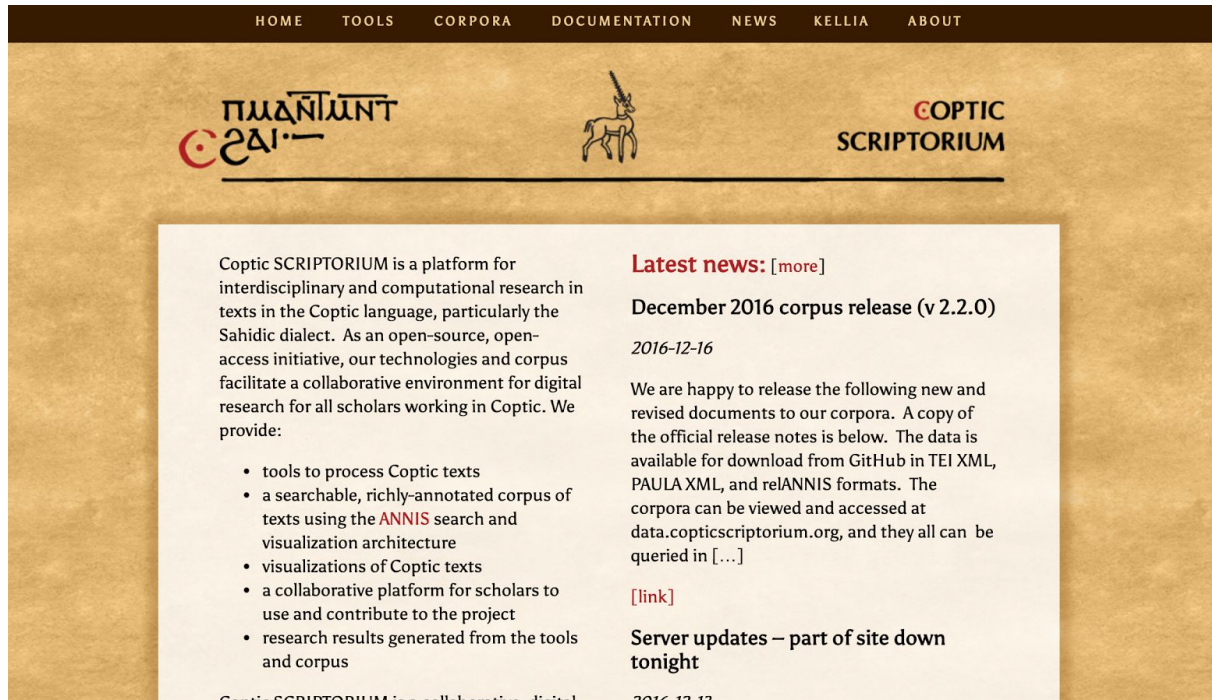
Alain Delattre, Assistant Professor, Department of Languages and Literatures, Université libre de Bruxelles; Papryi.info.

Eitan Grossman, Assistant Professor, Department of Linguistics and the School of Language Sciences, Hebrew University.

Robin Imhof, Humanities Librarian and Associate Professor, University Library, the University of the Pacific.

Appendix B: Website Screenshots

Main website and portal at www.copticscriptorium.org:



Web application for citation and access using CTS URNs at data.copticscriptorium.org:



Coptic Scriptorium Corpora & URN Resolver

Coptic Scriptorium provides Coptic texts for reading, analysis, and complex searches. The texts are citable and accessible through stable URNs, such as `urn:cts:copticLit:shenoute.fox` for Shenoute's work *Not Because a Fox Barks*. This application will provide the most recent version of our documents in the formats currently available for each text. If you know the URN for the material you seek, please enter it in the box above. You may also find documents by using the filters on the menu above.

If you wish to read Coptic texts, you can view individual documents online (in HTML) in various visualizations. We also provide links to our corpora in our [search tool ANNIS](#), as well as data files in [TEI XML](#), [PAULA XML](#), and [ANNIS](#) formats. This web application provides the most recent version of the data. Previous versions are available on [GitHub](#).

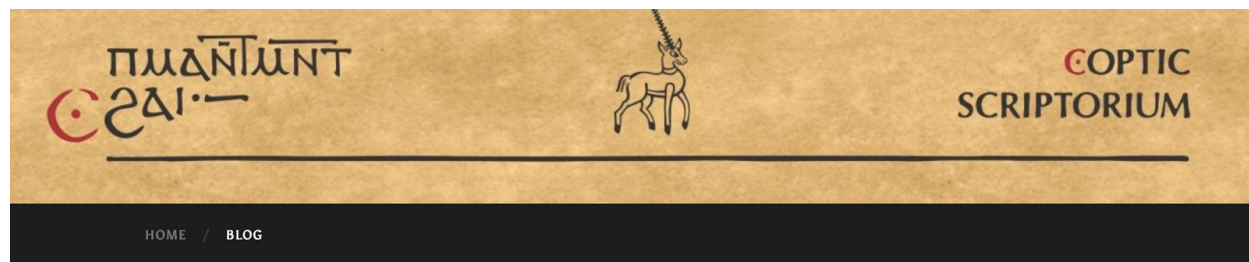
Documentation beneath the text provides information about it, including the version and date of publication online.

Resources

[Documentation on our part of speech tagging](#)

[Citation guidelines](#) for citing these documents

Project blog at blog.copticscriptorium.org:



December 2016 corpus release (v 2.2.0)

DECEMBER 16, 2016 / CTSCHROEDER / 0 COMMENTS

We are happy to release the following new and revised documents to our corpora. A copy of the official release notes is below. The data is available for download [from GitHub](#) in TEI XML, PAULA XML, and reANNIS formats. The corpora can be viewed and accessed at data.copticscriptorium.org, and they all can be queried in [ANNIS](#). We plan for another release with more documents in

blog.copticscriptorium.org/2016/12/16/december-2016-corpus-release-v-2-2-0/

Search form

SEARCH

RECENT POSTS

[December 2016 corpus release \(v 2.2.0\)](#)

[Server updates – part of site down tonight](#)

[NEH White Paper \(Preservations and Access Grant\) published](#)

Project Wiki at wiki.copticscriptorium.org:

ANNIS search and visualization interface landing page:

The landing page includes a search bar on the left with the placeholder text "Please enter AQL query". Below it are buttons for "Search", "More", and "History". A welcome message states: "Welcome to ANNIS! A tutorial is available on the right side." Below this is a "Corpus List" section with a "Visible" dropdown set to "scriptorium". A table lists various corpora with columns for Name, Texts, Tokens, and icons for information and download.

Name	Texts	Tokens	Info	Download
apophthegmata.patrum	36	4,217	Info	Download
besa.letters	3	2,296	Info	Download
coptic.treebank	10	4,361	Info	Download
doc.papyri	3	290	Info	Download
sahidica.1.corinthians	16	12,471	Info	Download
sahidica.mark	16	20,185	Info	Download
sahidica.nt	259	233,906	Info	Download
shenoute.a22	4	7,589	Info	Download
shenoute.abraham.our.father	7	7,670	Info	Download
shenoute.eagerness	10	10,030	Info	Download
shenoute.fox	1	2,814	Info	Download

The center section displays "Example Queries" with a table of query examples, their descriptions, and links to open corpus browsers.

Example Query	Description	open corpus browser
<code>Q norm="noyṛe"</code>	search for the normalized word noyṛe	shenoute.a22
<code>Q pos="NPROP"</code>	search for proper names	shenoute.abraham.our.father
<code>Q figure</code>	search for figures in the papyrus	doc.papyri
<code>Q gap_reason="illegible"</code>	search for illegible gaps in the text	doc.papyri
<code>Q norm="ana"</code>	search for the normalized word ana	apophthegmata.patrum
<code>Q norm="tonoc"</code>	search for the normalized word tonoc	shenoute.eagerness
<code>Q morph="peq"</code>	search for words containing the morpheme peq	shenoute.eagerness
<code>Q pos="FM"</code>	search for words tagged as foreign material	shenoute.eagerness
<code>Q norm="noyṛe"</code>	search for the normalized word noyṛe	shenoute.abraham.our.father
<code>Q pos="NPROP"</code>	search for proper names	shenoute.a22
<code>Q pos="VSTAT"</code>	search for stative verb forms	apophthegmata.patrum
<code>Q lang="Greek"</code>	search for words of Greek origin	shenoute.a22
<code>Q norm="anok"</code>	search for the unit "anok"	sahidica.nt
<code>Q norm="cofia"</code>	search for the normalized unit "cofia"	sahidica.nt
<code>Q lang="Greek"</code>	search for automatically detected Greek lexemes	sahidica.nt
<code>Q norm="nkotk"</code>	search for the normalized word nkotk	besa.letters

ANNIS search and visualization tool query results:

The results page shows the query `norm="noyṛe"` in the search bar. The results are displayed in a list of matches, with the first match selected. The detailed view of the first match shows the text "μηλοεις · εἰρηνοῦτε ὁντ ἐγνος ἡμῶν" and provides options for annotations, analytic view, diplomatic text, and normalized text. The list of documents shows 52 matches in 4 documents, with the first document selected.

Name	Texts	Tokens	Info	Download
sahidica.mark	16	20,185	Info	Download
sahidica.nt	259	233,906	Info	Download
shenoute.a22	4	7,589	Info	Download
shenoute.abraham.our.father	7	7,670	Info	Download
shenoute.eagerness	10	10,030	Info	Download
shenoute.fox	1	2,814	Info	Download

Appendix C: Data Visualizations

Normalized visualization

Normalized Text

[XH204] ερωαντβαωρ αωωκακ εβολ αν ετεντοκ πε πζηαλ μπμμμωνας
ζηζηηηροογ εγωω , ερεπμογιτρε ετεανοκ πε πζηαλ μπεχριστοσ , †σοογν
δεεκ†ογβηι αν , αλλα εκ†ογβειησογσ ετογνηζ
ρωωε εροκ ντοκ μππεκειωτ πδιαβολοσ ετογν
μππεγειωτ ετογνηζ νζητογ ιησογσ ετογκω νζητογ
μμμγ ωπιε ενεζ ντρε ντακχοοσ , πωπιε μπειμα , πεοογ μπειμα αιπαραιτει μμογ ,
†σοογν γαρ επετισωωτ εβολ ζηητ . νετεογνταγπεοογ μνιπταειο εβολ ζιτνιησογσ .
εγρογ νεοογ ζιταειο νρωμε . ντρε γαρ ετεμνμντληστησ ωοοπ ννετεογνταγ ιησογσ
ζηνογμε καταπεντακχοογ εροι εβολ δεαιφι ννεκνογτε ζηνογσραστ αγω δεαιτρεγμογρ
μπεκσωω μππεκωπιε εζογν ενογεσρο μπεκνι εγσζη εζενχαρτησ · εαγογωσπ
ννεκμνμμογ ετζηνηεωωωγ ζωσ νρπ ζιχμμπνι μπεκνι . αγω εζογν ζμπεκρο , μνπρο
ννετεινε μμοκ , εμμνμντρμζε ωοοπ ννετκω νζηηγ εκρονοσ , ετε [XH206] ντοκπε
μννεττντων εροκ ζημμνταπιστοσ , μνμντακαθαρτοσ νιμ . ογ νζωβ η αω νωαδε
ζηννεντακχοογ νετο μμντρε εροκ αν δεεκνι επσατανασ . νεπιστολη νε εντακπαζογ
 , η †σοογν αν μπατιχοογσογ δεκναπαζογ , πεωβρρζωβ ννενταγωωωτ ζμπτοκ
μπεγραμματαεγσ ννσελισ μπλχωωμε ννωαδε ντανεπροφητησ χοογσε ναγ ζμπραν

It's not when the fox cries out, which is you,
oh servant of Mammon, in voices that shout,
that the lion, which is I, the servant of Christ,
is afraid.

Shenoute, *Not Because a Fox Barks*, ed. Amir Zeldes and Caroline T. Schroeder. Trans. Amir Zeldes. *Coptic SCRIPTORIUM*. urn:cts:copticLit:shenoute.fox.monbxh_204_216. v. 1.5, 13 May 2106. http://data.copticscriptorium.org/urn:cts:copticLit:shenoute.fox.monbxh_204_216.

Diplomatic Visualization

Diplomatic Edition

XH205

	15 ԵՐԾԱՆԿԵՎՈՐ՝ ապկաքեբոլ՝ ան՝ետե՛նտոկ պե՛րբալ միմամոնա՛ս
	20 ջնջեղբոսոյ եղօօ՝,բրե պիօյի՛րբե՛ ետեանօկպե՛ րբալմի
	25 ԽԵ,ԺՏՕՕԿՆ չեքԺօյնի ան՝ալաեկԺ օյբեւեօյ նջնեքի
ՏԻԱՆՕ՝.ԻՇ	Ռեգարեմն
ՕՆ՝րօյբերօկ	մնտլնստի՛ս
նտօկմիք	օթօթնեթե
եւտլլաւօ	օյնտալի՛սնօյ
5 ԼՕ՝ետօյնջ՝	5 Ե՛կաթաթե
նջնի՛կե՛տ	տալ.օթի՛րօի
ջեւնչե՛րօյ՝	եւօլ.չեւի՛ն
նտօյջօյ	նեկնօյթե
միպեւօտ՝	ջնօյսրա՛շտ
10 ԵՏՕԿՆջն	10 ԼԿալալթրեյ
տօյի՛ետօյկա՝	նօյր՝նեք
նջտիյերօյ՝	սօյմիք
միմտեօնալ	սիպեթօն՝

Shenoute, *Not Because a Fox Barks*. *Coptic SCRIPTORIUM*.

urn:cts:copticLit:shenoute.fox.monbxh_204_216.

Analytic Visualization

ACOND ART N V ADV NEG CREL PPERI COP ART N PREP ART
 ερωαν| τ| βαυορ αωακακ εβολ αν ετε| ντοκ πε π| ζμζαλ μ| π|
NPROP PREP ART N CREL PPERS VSTAT PUNCT CLOC ART N V CREL PPERI COP
 μαμμωνας ζν| ζεν| ζροογ ε| γ| ου , ερε| π| μογι| τρε ετε| ανοκ πε
ART N PREP ART N PUNCT
 π| ζμζαλ μ| πε| χριστος ,

It's not when the fox cries out, which is you, oh servant of Mammon, in voices that shout, that the lion, which is I, the servant of Christ, is afraid.

PPERS V CONJ CLOC PPERS V PREP PPERO NEG PUNCT CONJ CLOC PPERS V PREP NPROP CREL
 †| σοογν δε| ε| κ| †| ουβη| ι αν , αλλα ε| κ| †| ουβε| ιχογς ετ
VSTAT PREP ART N PUNCT
 | ουηζ ζν| νε| χρειστιανος .

I know it's not against me you fight, but against Jesus who dwells inside the Christians.

NPROP ADV V PREP PPERO PPERI PREP PPOS N ART N CREL VSTAT PREP
 ιχογς ον ρωθε ερο| κ ντοκ μν| πεκ| ειωτ π| διαβολος ετ| ουηζ νζητ|
PPERO CREL PPERS V PREP PPERO PUNCT
 κ ετ| κ| ζελπιζε ερο| ς .

Jesus still conquers you yourself and your father the devil who dwells inside you, for which you hope.

PPERI IMOD PPERO PREP PPOS N CREL VSTAT PREP PPERO N CREL PPERS V PREP N
 ντοογ ζω| ου μν| πεγ| ειωτ ετ| ουηζ νζητ| ου ιχογς ετ| ου| κω ν| ζτη
PPERO PREP PPERO PUNCT
 | γ ερο| ς .

And they too on their part, with their father who dwells inside them, Jesus, on whom they depend.

Shenoute, *Not Because a Fox Barks*. *Coptic SCRIPTORIUM*.

urn:cts:copticLit:shenoute.fox.monbxh_204_216.

Chapter visualization

Sahidica Chapter View

1: ἀνοκ δε ρω ντερει ωαρωτν νασνηγ νταιει ρνογχιςε αν νωαχε η νκοφια ειχω
ερωτν ντμντμντρε μπινογτε .

2: μπιμееγe γαρ χετσοογν νλααγ νε εαγσταγρογ μμογ .

When I came to you, brothers, I didn't come with excellence of speech or of wisdom, proclaiming to you the testimony of God.

3: ἀνοκ ρω νταιει ωαρωτν ρν ογμντςωβ . μνογρστε . μνογςτωτ εναωωγ .

4: αγω παωαχε μνπαταωεοειω νταφωωπε αν ρνογπιθε νκοφια νωαχε . αλλα ρνογογωνρ εβολ μπινεγμα ριςομ .

5: δεκαας εννετνπιςτις ωωπε ρνογπειθε νκοφια νρωμε αλλα ρνογςομ ντεπινογτε .

6: ενωαχε δε νογςοφια ρνντελειος ογςοφια δε ενταπειαιων αν τε ογδε νταναρχων αν τε μπειαιων ναι ετναογωςγ .

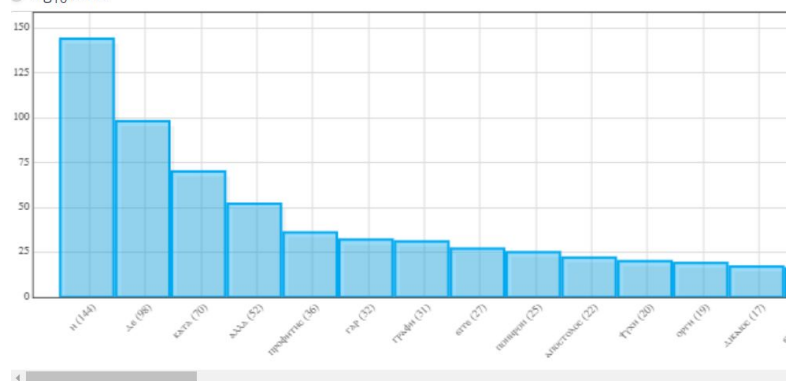
7: αλλα ενωαχε νογςοφια ντεπινογτε ρνογμγςτηριον . ται ετρηπ τενταπινογτε πορς εβολ ρατρη νηαιων επενεοογ .

1 Corinthians 2 (Sahidica version), ed. Elizabeth Platte et al., taken from J. Warren Wells' Sahidica edition. Trans. World English Bible. *Coptic SCRIPTORIUM*.
urn:cts:copticLit:nt.1cor.sahidica_ed:2. v. 1.1, 22 May 2105.
http://data.copticscriptorium.org/urn:cts:copticLit:nt.1cor.sahidica_ed:2.

Histogram visualization for frequency analysis from a query in the ANNIS database

Showing histogram of top 100 results, see table below for complete dataset.

☒ linear scale
☐ log₁₀ scale



289 items with a total sum of 1337 (query on shenoute.eagerness, shenoute.fox, shenoute.abraham.our.father, sh

rank	#1 lemma	count
1	η	144
2	Δε	98

Appendix D: Digital Coptic 2 Program

Digital Coptic 2 Program

Thursday, March 12: Symposium

9 am-6 pm

Georgetown University

Poulton Hall, Room 230

1421 37th St. NW

Washington, DC 20057

Presentations are each 20 minutes long followed by 10 minutes for questions and discussion.

9:00 Coffee and Arrivals

9:15 Welcome

Amir Zeldes, Georgetown University

9:30 New and Expanding Digital Projects in Coptic Studies

Chair: Caroline T. Schroeder, the University of the Pacific

Adventures in Crowd-Sourcing Papyri - The Resurrecting Early Christian Lives DH Project, Philip Sellev, the University of Minnesota

Website Galleries of the White Monastery Candle Room Manuscript Fragments: Challenges of Digitization and Classification, Mary K. Farag, Yale University

The Digital Edition of the Coptic-Sahidic Old Testament and its planned Virtual Manuscript Room (VMR), Frank Feder, Akademie der Wissenschaften zu Göttingen

Digitizing Language Contact: Lexicography and Technological Perspectives at the Database and Dictionary of Greek Loanwords in Coptic (DDGLC), Frederic Krueger and Katrin John, Universität Leipzig

Discussion (30 minutes)

12:00-1:00 Lunch

1:15 Digital and Computational Research in Coptic Language and Literature

Chair: Elizabeth Platte, Valparaiso University

Coptic SCRIPTORIUM: Current Possibilities and Future Directions, Amir Zeldes, Georgetown University, and Rebecca S. Krawiec, Canisius College

Coptic Scriptorium beyond the Manuscript: Tokenization and Corpus Analysis, Paul Dilley, the University of Iowa

Synthesis, Boundness, and Clitics in Sahidic Coptic, So Miyagawa, Kyoto University

Discussion (30 minutes)

3:15 Coffee Break

3:30 Digital Humanities and Eastern Christian Traditions beyond Coptic Studies

Chair: Christine Luckritz Marquis, Union Presbyterian Seminary

Digital Preservation and Oral History of Displaced Syriac Speakers in the Middle East, Robin Darling Young, the Catholic University of America

A New XML Exchange Format for Aligning Translations, Quotations, and Other Versions of Texts, Joel Kalvesmaki, Dumbarton Oaks

Ex uno pro pluribus: Digitization, cataloging, and study of Eastern Christian manuscript collections at the Hill Museum & Manuscript Library, Adam Carter McCollum, Hill Museum and Manuscript Library

Discussion (30 minutes)

5:30 Concluding Remarks and Wrap-Up

Friday March 13

9 am-6 pm

Georgetown University

Poulton Hall, Room 230

1421 37th St. NW

Washington, DC 20057

A day-long workshop on Coptic SCRIPTORIUM for the SCRIPTORIUM team, collaborators, contributors, and those interested in becoming collaborators or contributors. We will discuss SCRIPTORIUM technologies, how to contribute and annotate text corpora for the project, future directions, and possible collaborations. A final agenda will be distributed in March. Space is limited. Please indicate on the registration form if you wish to join us, and we will confirm your attendance.

Appendix E: Guidelines

The following documents are attached:

- Citation Guidelines
- Transcription Guidelines
- Part of Speech Tagging
- Lemmatization Guidelines

Citation Guidelines

When researching our corpora for a future publication, **please note the date and version number of the documents or corpora while you are conducting your research.** *We update our corpora regularly* and recommend all citations include the version number and date of the resources used, as described below. (If you conducted research in the past and did not note the version and date of the corpus at that time, then please cite the date you accessed the corpus.)

Examples from the Chicago Manual of Style format appear below; please modify as appropriate for other style formats.

Cite this project

We recommend upon first citation of the project, text, query, tool, or other resource, you cite the full project as well as the individual resource.

First citation

Caroline T. Schroeder, Amir Zeldes, et al., *Coptic SCRIPTORIUM*, 2013-[current year], <http://copticcriptorium.org>.

Bibliography

Schroeder, Caroline T., Amir Zeldes, et al., *Coptic SCRIPTORIUM*. 2013-[current year]. <http://copticcriptorium.org>.

Cite a corpus or corpora

First citation

Coptic SCRIPTORIUM, [*corpus name*], [corpus URN], [version number], [date]. [http://data.copticcriptorium.org/\[corpus URN\]](http://data.copticcriptorium.org/[corpus URN]).

E.g.,

Coptic SCRIPTORIUM, *apophthegmata.patrum Corpus*, urn:cts:copticLit:ap, v. 1.5, 4 October 2015. <http://data.copticscriptorium.org/urn:cts:copticLit:ap>.

Subsequent citations

Coptic SCRIPTORIUM, [corpus urn].

E.g.,

Coptic SCRIPTORIUM, urn:cts:copticLit:ap.

Bibliography

Coptic SCRIPTORIUM. [*corpus name*]. [corpus URN]. [version number], [date].
[http://data.copticscriptorium.org/\[corpus URN\]](http://data.copticscriptorium.org/[corpus URN]).

E.g.,

Coptic SCRIPTORIUM. *apophthegmata.patrum Corpus*. urn:cts:copticLit:ap. v. 1.5, 4 October 2015. <http://data.copticscriptorium.org/urn:cts:copticLit:ap>.

Cite an individual document

First citation

Author, *Ancient title* [chapter.verse if available], ed. [annotators], trans. [translation].
Coptic *SCRIPTORIUM*. [document URN]. [version number], [date].
[http://data.copticscriptorium.org/\[document URN\]](http://data.copticscriptorium.org/[document URN]).

E.g.,

Sahidic Apophthegmata Patrum 6, ed. Paul Lufter et al., trans. Paul Lufter and Amir Zeldes. *Coptic SCRIPTORIUM*. urn:cts:copticLit:ap.6. v. 1.1.0, 21 May 2015.
<http://data.copticscriptorium.org/urn:cts:copticLit:ap.6>.

Besa, *Letter to Aphthonia* 1.1, ed. Amir Zeldes and Caroline T. Schroeder, trans. Amir Zeldes. *Coptic SCRIPTORIUM*. urn:cts:copticLit:besa.aphthonia. v. 1.3.0, 28 May 2015.
<http://data.copticscriptorium.org/urn:cts:copticLit:besa.aphthonia>

Subsequent citations

Author, *Ancient title (abbreviated)* [chapter.verse if available], *Coptic SCRIPTORIUM*, [document urn].

E.g.,

Sahidic AP 6, *Coptic SCRIPTORIUM*, urn:cts:copticLit:ap.6.

Besa, *Aphthonia* 1.1, *Coptic SCRIPTORIUM*, urn:cts:copticLit:besa.aphthonia.

Bibliography

Author. *Ancient title*. Ed. [annotators]. Trans. [translation]. *Coptic SCRIPTORIUM*. [document urn]. [version number], [date]. [http://data.copticscriptorium.org/\[document URN\]](http://data.copticscriptorium.org/[document URN]).

E.g.,

Sahidic Apophthegmata Patrum. Ed. Paul Lufter et al. Trans. Paul Lufter and Amir Zeldes. *Coptic SCRIPTORIUM*. urn:cts:copticLit:ap.6. v. 1.1.0, 21 May 2015.
<http://data.copticscriptorium.org/urn:cts:copticLit:ap.6>.

Besa, *Letter to Aphthonia*. Ed. Amir Zeldes and Caroline T. Schroeder. Trans. Amir Zeldes. *Coptic SCRIPTORIUM*. urn:cts:copticLit:besa.aphthonia. v. 1.3.0, 28 May 2015.
<http://data.copticscriptorium.org/urn:cts:copticLit:besa.aphthonia>

Citing and linking to a query

Links to queries in ANNIS are possible; however, these links remain stable only until the corpus is revised or updated. We recommend **citing the project** as well as providing information about **the query, query link, and date and/or version accessed**. We also recommend using the download option in ANNIS to **download query results**.

Researchers may then save (or even publish on their own websites under the proper license) the raw data for future reference.

Coptic SCRIPTORIUM Diplomatic Transcription Guidelines

Version: 1.2_2016.8.26

Caroline T. Schroeder¹ & Amir Zeldes²

1. University of the Pacific

2. Georgetown University

1. Preamble

This document details guidelines for transcribing a diplomatic edition of a manuscript in Sahidic Coptic according to the Coptic SCRIPTORIUM project scheme. The diplomatic transcription currently requires extensive manual annotation, due to the complexities of processing a diplomatic text in which no word breaks exist in the original and yet words and even morphemes span across line, column, and page breaks.

The transcription procedure assumes familiarity with basic paleography and traditional manuscript transcription following the Leiden conventions.

(<http://www.stoa.org/epidoc/gl/latest/app-glossary.html#leiden>)

The diplomatic transcription also utilizes XML (eXtensible Markup Language) -like tagsets, including some of the TEI (Text Encoding Initiative) XML markup language, although the resulting document is **not** a valid XML document. Wherever possible, the EpiDoc subset of TEI XML is utilized for element nomenclature. EpiDoc TEI conventions were created by and for epigraphers and have come to be a standard in markup of ancient texts, epigraphic or otherwise.

(<http://sourceforge.net/p/epidoc/wiki/Home/>) In contrast to TEI, SCRIPTORIUM utilizes no milestone XML tags (e.g., <cb/>). Instead, all tags are span annotations (e.g., <cb>This is a column of Coptic text.</cb>).

We recommend using an XML editor such as Oxygen to ensure the encoding is well-formed and well-structured.

The aim is twofold: 1) to achieve a transcription that documents the text and visualization of the manuscript as closely as possible to the original; 2) to provide a text file that can be processed by various digital tools and software, such as a tokenizer, a part-of-speech tagger, or the ANNIS database infrastructure (<http://www.sfb632.uni-potsdam.de/annis/>; Zeldes et al. 2009). Coptic SCRIPTORIUM has bundled some of these tools in a [Natural Language Processing web service](#).

The resulting transcription itself does not resemble a traditional text of a diplomatic edition. The markup ensures optimization for processing and search using such tools and software. For examples of the diplomatic editions visualized in HTML generated from the post-ANNIS transformations, see corpora at data.copticscriptorium.org. Valid EpiDoc TEI XML versions of the documents are also provided from this site.

2. Character Encoding

Texts are encoded using the UTF-8 (Unicode) Coptic language character set. The freely available Antinoou font and Coptic-English keyboard created by Michael Everson in cooperation with the International Association of Coptic Studies is the standard (<http://www.evertype.com/fonts/coptic/>). Unicode characters in the private use area are not recommended.

2.1 Alphanumeric Characters

Characters follow the orthography of the manuscript.

Mark oversize characters with XML tagging. Do not use uppercase version of the character.

2.2 Punctuation and Decoration

Punctuation and decoration follows the manuscript as closely as possible within the Unicode character set. Not all decoration and punctuation can be encoded using characters; deviations or documentation that can't be keyed in is instead typically indicated in a note element.

Notes on individual specific punctuation characters:

For the character ` that occasionally appears at the end of words in some manuscripts, use U+2CFF. Example:

ⲡⲉⲙⲙⲟⲛ`ⲧⲉ

2.3 Accentuation and Supralinear Strokes

Accentuation and supralinear strokes follow the orthography of the manuscript. Some manuscripts have binding strokes between letters (e.g. ⲉⲛ̅) whereas others in the case of the same word might only provide a stroke over a single letter (e.g., ⲉⲛ̅). The diplomatic transcription follows the conventions of the manuscript, even if the manuscript is internally inconsistent or contains what seem to be errors.

Notes on encoding individual specific accents, strokes, etc, using the Coptic-English keyboard for Antinoou (for MacIntosh):

- (as in ⲛ̅) the supralinear stroke above only one letter: type the letter followed by Unicode U+0304 (; on keyboard)
- (as in ⲛ̅ⲛ̅) the binding stroke between two letters: type first letter then U+FE24 (< in the Coptic-English keyboard) then second letter then U+FE25 (> in the Coptic-English keyboard), i.e. m<n> on a Mac using the Coptic-English keyboard
- (as in ⲛ̅ⲛ̅ⲧ̅) binding stroke over three letters: type the first letter then U+FE24 (< on a Mac using the Coptic-English keyboard) then second letter then U+FE26 (: [i.e. shift+;] on a Mac using the Coptic-English keyboard) then third letter then U+FE25 (> on a Mmac using the Coptic-English keyboard), i.e. m<n:t>

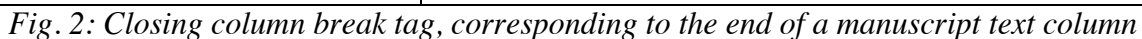
For squiggly curved or jagged strokes over etas, use a regular circumflex rather than a dot or line or trema (ˆ): type the letter followed by U+0302 (option+3 on the keyboard)

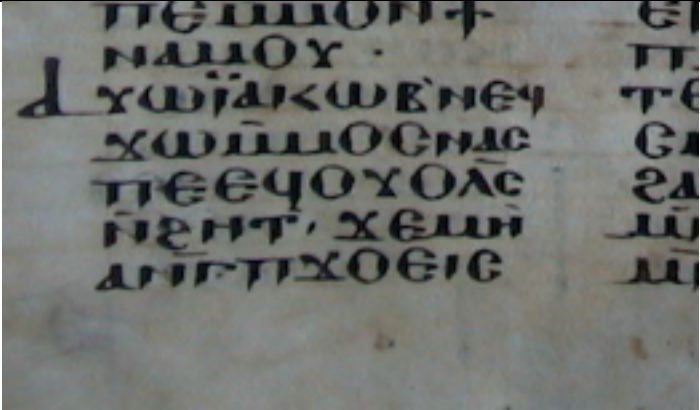
3. Text Divisions

3.1 Line Breaks

3.2 Column Breaks

Fig. 1: Opening column break tag, corresponding to beginning of manuscript column



<p> <code> ⲁⲛⲣⲉ̀_ⲉ ⲁⲱ ⲡⲉ_ⲙⲏⲟ ⲛ`_ⲧ ⲛⲁ ⲛⲟϥ_·_ ⲉⲗⲱ_ⲓⲁⲕⲱⲃ`_ⲛⲉ ⲥ ⲕⲱ_ⲙⲏⲟ ⲥ_ⲛⲁ ⲥ_ ⲡⲉ_ⲉ ⲥ ⲟϥⲟⲗⲥ_ ⲛ ⲣⲏⲧ_·_ⲕⲉ_ⲙⲏ`_ ⲁⲛⲧ ⲡ ⲕⲟⲉⲓⲥ_</cb> </code> </p>	
<p>MONB.YA 520; Coptic Manuscript IB 2 f. 27v, Naples, Biblioteca Vittorio Emanuele III,”</p>	

3.3 Page Breaks and Numbering

All page breaks in the transcription should follow the page divisions in the manuscript.

Page numbering in the transcription reflects the page numbering in the original manuscript codex. Codex sigla in the example below are two-letter codes following the White Monastery codex siglum list created by Tito Orlandi (Orlandi 2002; also <http://www.cmcl.it/>). Page breaks are wrapped in TEI compatible span annotations using the `<pb></pb>` tagset with the `xml:id` element. The entire page of text (including the relevant column tags) should be wrapped with these tags. Thus `<pb xml:id="YA518">` is the opening tag for page 518 in White Monastery codex YA (MONB.YA). The `xml:id` should not contain spaces. (Thus, `xml:id="YA518"` not `xml:id="YA 518">`.)

Fig. 3: Closing and opening page break tags indicating the end of one page and beginning of the next. (Note: the opening tag for the first page and closing tag for the second page are not visible here but are required.)

<p> <code> ⲑⲛⲟⲥ`_ⲉ ⲛⲁ ⲁⲱ ⲥ_·_ ⲉⲗⲱ_ⲡ ⲉⲧ ⲛⲏϥ_ ⲉⲃⲟⲗ_ⲉⲛ_ⲛⲉⲕ`_·_</cb></note></pb> <pb xml:id="YA519"><note note="page number ⲡⲓⲟ barely visible in upper right"><cb>[....]ⲥ_·_ⲓⲥⲁⲕ` [...]ⲥ_ⲡ ⲉⲧ ⲛⲁ [ⲕⲗⲏ]ⲣⲟⲛⲟⲙⲉⲓ_ ⲙⲏⲟ ⲕ_·_ ⲉⲧⲃⲉ_ⲡⲁⲓ_ⲣⲱ_ⲁ ⲥ ⲕⲟⲟ ⲥ_ⲉⲛ ⲟϥ </code> </p>
--

The location and Coptic numeration of the page number is currently documented in a note element. (See Figure 3 above).

4. Word Segmentation, Spacing, and Tokenization

Sahidic Coptic bound groups are formed by several words and/or morphemes attaching together. A word refers to one noun, preposition, article, etc. One complex word can be comprised of multiple morphemes, including affixes such as ⲁⲧ, ⲙⲛⲧ, or ⲡⲉⲓ, or compound words, such as complex numbers (e.g. -teens) and verbs formed with ⲡ. One bound group may include multiple prepositions and objects, or a verbal auxiliary + subject + infinitive, or even more words and morphemes strung together (generally speaking clitics). The copula, which some might consider a clitic, remains unbound. Coptic SCRIPTORIUM follows the practices in Bentley Layton's grammar (Layton 2011) for word, morpheme, and bound group segmentation.

Examples of individual words comprised of one morpheme:

Ⲫⲱⲧⲙ

ⲛⲟⲃⲉ

ⲉⲛⲧ

Examples of individual words comprised of multiple morphemes:

ⲙⲛⲧⲁⲧⲪⲱⲧⲙ

ⲡⲉⲓⲡⲛⲟⲃⲉ

Examples of bound groups comprised of words with multiple morphemes:

ⲧⲙⲛⲧⲁⲧⲪⲱⲧⲙ

ⲙⲡⲉⲓⲡⲛⲟⲃⲉ

ⲡⲣⲙⲛⲉⲛⲧ

ⲁⲛⲧⲁⲩⲧⲥ

Note: if a project wishes to annotate on the morpheme level (i.e. internal analysis of units like ⲙⲛⲧ) and not just on the *word* level, the morphemes need to be tokenized. Coptic SCRIPTORIUM annotates on the word level and then provides additional annotation on the morpheme level for compound words and words with affixes. (See section 4.4 for more information.)

In most manuscripts, no spaces between words or bound groups are provided. Sometimes a diacritical mark, such as ` does appear, but word segmentation following diacritics and punctuation does not always correspond with contemporary segmentation practices (such as Layton or Till (1960)). More study of this marking is required.

4.1 Word Segmentation

SCRIPTORIUM diplomatic transcription marks word segmentations according to Layton's conventions (Layton 2011). The transcriber inserts a unique character, such as an underscore ("_"), after each Coptic bound group, even when the end of the bound group falls at the end of a line.

Likewise, all punctuation is followed by an underscore.

- (1) ⲉⲧⲉⲓⲥⲙⲁⲛⲗ_ (word ends at end of line)
- (2) ⲡⲉ_ⲛⲓⲣⲏⲁⲕⲗⲏ (two words, in which the second bound group flows into line 3)
- (3) ⲣⲟⲛⲟⲙⲉⲓ_ⲙ (the bound group continues from line 2, is followed by an underscore)
- (4) ⲙⲟⲕ_ⲁⲛ_·_ (punctuation followed by an underscore)

These underscores are not and do not need to be visualized in HTML transformations of the diplomatic editions; they are nonetheless essential for processing the text, since they demark breaks between bound groups and will enable searches and visualizations of a word-segmented text.

We do not recommend using spaces to demarcate bound groups and punctuation, since spaces may occur elsewhere in the document (such as inside XML tags), and lead to confusion during automatic processing.

4.2 Spacing

Encoding of blank space is preferred to using the space key. The encoding should match spaces in the manuscript. Consequently, if the manuscript provides no spaces between words or punctuation, the diplomatic transcription contains no spaces. Where there are significant spaces in the manuscript that the transcriber wishes to draw attention to, the transcription should encode a space using TEI XML tags in order to visualize the white space in the manuscript. Encode the word, morpheme, or punctuation next to the white space, as in these examples:

- (1) <hi rend="1_space_right">·</hi> will visualize one space to the right of the ·
- (2) ⲛ<hi rend="1_space_right">ⲃⲟⲛⲏⲟⲥ</hi> will visualize one space to the right of he n t
- (3) ⲁⲱ<hi rend="2_space_right">ⲓ</hi>_<hi rend="1_space_right">·</hi>_ⲉⲧⲃⲉ_ will visualize two spaces to the right of ⲓ and one space to the right of the ·

It is important to make sure that attributes are surrounded by straight, not curly quotes (i.e. " on both sides).

4.3 Tokenization of Words

If one wishes to manually segment bound groups into words, one can do so using the pipe character (“|”).

- (1) ⲕⲓⲧⲙ|ⲡ|ⲛⲟϥⲧⲉ_ (preposition|article|noun)
- (2) ⲉⲧⲉ|ⲓⲥⲙⲁⲛⲗ_ (converter|noun)
- (3) ⲡⲉ`_ⲛ|ⲓ|ⲛⲁ|ⲕⲗⲏ (word_auxiliary|subject pronoun|future marker|verb (verb continues to line 4))
- (4) ⲣⲟⲛⲟⲙⲉⲓ`_ⲙ

The NLP web service contains a tokenizer that will take as input bound groups and provide as output word segmentation with pipes. Coptic SCRIPTORIUM’s standalone tokenizer tool will do the same.

4.4 Tokenizing and Annotating Morphemes below the Word Level

To conduct research on the morpheme level in compound words or other words that contain multiple morphemes, the words will need to be tokenized and annotated below

the word level and on the morpheme level. In Coptic SCRIPTORIUM, text is annotated on the word level for the part of speech (see [SCRIPTORIUM Part-of-Speech Tagsets for Sahidic Coptic](#)) and other characteristics, such as language of origin. Tokenizing and annotating on the morpheme level allows for additional search, visualization, and research capabilities.

Examples of individual words comprised of multiple morphemes, tokenized on the morpheme level:

word	ⲙⲛⲧⲁⲧϥⲱⲧⲙ		
morpheme	ⲙⲛⲧ	ⲁⲧ	ϥⲱⲧⲙ

word	ⲣⲉϥⲣⲛⲟⲃⲉ		
morpheme	ⲣⲉϥ	ⲣ	ⲛⲟⲃⲉ

Examples of bound groups comprised of words with multiple morphemes:

bound group	ⲧⲙⲛⲧⲁⲧϥⲱⲧⲙ			
word	ⲧ	ⲙⲛⲧⲁⲧϥⲱⲧⲙ		
morpheme	ⲧ	ⲙⲛⲧ	ⲁⲧ	ϥⲱⲧⲙ

bound group	ⲙⲡⲣⲉϥⲣⲛⲟⲃⲉ				
word	ⲙ	ⲡ	ⲣⲉϥⲣⲛⲟⲃⲉ		
morpheme	ⲙ	ⲡ	ⲣⲉϥ	ⲣ	ⲛⲟⲃⲉ

bound group	ⲡⲣⲙⲛⲉⲛⲧ			
word	ⲡ	ⲣⲙⲛⲉⲛⲧ		
morpheme	ⲡ	ⲣⲙ	ⲛ	ⲉⲛⲧ

Compound words that involve an article or affixed personal pronoun to the second item of the compound typically are tokenized as bound groups comprised of multiple words, not as one word comprised of multiple morphemes.

Examples of bound groups containing compound words with articles or pronouns on the second unit of the compound:

bound group/compound	PḲNΔϣ		
word	P	ḲNΔ	ϣ
<i>no tokenization & annotation on the morpheme level below the word level</i>			

bound group	MΠETNPTMEEY				
word	MΠE	TN	P	Π	MEEY
<i>no tokenization & annotation on the morpheme level below the word level</i>					

(where PTMEEY is considered to contain multiple words, not morphemes below one word level)

[Note: the part-of-speech tagger developed by Coptic SCRIPTORIUM operates on the *word* level, not the sub-word morpheme level. So, PḲOT is tagged as one V, MNTATCOTM as one N, etc.]

Transcription conventions for segmenting morphs should utilize a unique character, such as a dash or hyphen. E.g.:

T|MNT-ΔT-COTM

M|Π|PEY-P-NOVE

If you plan to use Coptic SCRIPTORIUM's NLP web service, you may transcribe the Coptic in bound groups with no pipes or morphemes. The NLP web service's tokenizer can provide as output segmentation with pipes between words and dashes between morphs. Likewise, Coptic SCRIPTORIUM's stand-alone tokenizer can output words with segmented morphs. The webservice can further automatically annotate the segmented words and morphs for part of speech, language of origin, and lemma.

5. Rendering and Leiden Transcription Conventions

Coptic SCRIPTORIUM uses Leiden and Leiden+ conventions for transcribing manuscripts. The encoding follows the EpiDoc guidelines. Not all Leiden documentation is currently XML encoded as Leiden+, however.

5.1 Characters Highlighted, Raised, Lowered, or Set Apart in Some Way

Characters that are raised, lowered, or printed in different colors or styles are encoded using the TEI XML element <hi> with the rend attribute. Letters written above the line are encoded: <hi rend="superscript">. Characters written below the line are encoded: <hi rend="subscript">. Letters in a different color ink are encoded with the color ink, e.g., <hi rend="red">. It is possible to combine these annotations, e.g. <hi rend="red subscript">. Coptic SCRIPTORIUM currently encodes large, tall (the letter stretches above the line), long (letter stretches below the line), thin, superscript, subscript, and colors. Any additional information can be provided in a note element. To encode two attributes, use a space (not a comma) between the two attributes.

Example

ⲉⲓⲧⲙⲡⲛ<hi rend="superscript"><note note="o is
directly above the γ">o</note></hi>γ
ⲡⲡⲉⲧⲛⲁⲛⲟ<hi rend="large">γ</hi><hi rend=
"long thin">ϥ</hi>_._

Diplomatic Visualization

ⲉⲓⲧⲙⲡⲛϥ (ANNIS) or
ⲉⲓⲧⲙⲡⲛ\o/γ (EpiDoc XSLT)

ⲡⲡⲉⲧⲛⲁⲛⲟϥ.

Other encodings are colors (red, brown, green, etc.) “ekthetic” should be used for characters that are part of the ongoing text but written to the left of the margin line. See below, in which the ⲡ is encoded <hi rend="red large ekthetic">ⲡ</hi>

ⲉⲧⲁⲛⲕⲉⲃⲟⲗⲓⲛⲧⲉ
ⲛⲡⲉⲃⲃⲓⲟⲛ ⲉⲛⲧⲓⲛ
ⲡⲁⲓⲁⲉⲁⲥⲁⲛⲟⲩⲱⲁ

hi@rend cannot contain more than five words as per Epidoc guidelines and may contain only alphanumeric characters. (No punctuation. So <hi rend= "long, thin">ϥ</hi> is invalid.)

5.2 Damaged Characters

Characters that are damaged but restored based on context are marked with an underdot. Coptic SCRIPTORIUM uses the diacritical character ̣ (Unicode U+0323). These characters are not currently encoded in TEI XML using the EpiDoc tagset for Leiden+. Coptic SCRIPTORIUM uses the underdot character rather than annotation to designate this information.

5.3 Lacunae and Lost Characters

Lost lines and characters (lacunae) are indicated using square brackets, as in the Leiden conventions. They may be encoded using the EpiDoc tagset, but it is not required. See EpiDoc guidelines for more details (“EpiDoc Guidelines: Lost Characters, Quantity Unknown”; “EpiDoc Guidelines: Editorial Restoration: Characters Lost but Restored by Modern Editor”; “EpiDoc Guidelines: Lost Characters, Quantity Approximate”; “EpiDoc Guidelines: Lost Characters, Quantity Known”; “EpiDoc Guidelines: Erased and Lost”; “EpiDoc Guidelines: Lacunas, Other Units”).

- (1) Example encoded using the gap element:

```
<gap reason="lost">
[ ]_
[ ]_
[ ]_
</gap>
```

- (2) Unencoded gaps (no XML elements):

```
[.....]ⲏⲣ[.]
```

5.4 Other

Other rendering information is encoded either according to EpiDoc conventions or recorded as information within a note element. See the cheatsheet for Leiden+ conventions in EpiDoc at <https://sourceforge.net/p/epidoc/code/HEAD/tree/trunk/guidelines/msword/cheatsheet.doc?format=raw> and http://papyri.info/docs/leiden_plus. See also the full list of text transcription guidelines here <http://www.stoa.org/epidoc/gl/latest/app-alltrans.html>.

Transcribing in Oxygen or a similar XML editor is recommended, to ensure tags are well-structured.

6.0 File Format and Document Preferences

Documents are transcribed in a text editor such as TextEdit. Document preferences are set to UTF-8 encoding without byte-order Mark (BOM). (E.g., in TextEdit 1.7.1 for Macintosh, in the File-->Preferences menu, click on “Open and Save,” and select “Unicode (UTF-8)” for Opening files and Saving files.)

Bibliography

An up-to-date bibliography can be found at the project’s Zotero page:

https://www.zotero.org/groups/coptic_SCRIPTORIUM/items/collectionKey/8IHTW3NZ

Bodard, Gabriel. “EpiDoc Appendix: Glossary: Leiden, Leiden-plus.” *Appendix: Glossary*. 18 Jun. 2013. <<http://www.stoa.org/epidoc/gl/latest/app-glossary.html#leiden>>.

“Corpus Dei Manoscripti Copti Letterari.” *CMCL - Studies in Coptic Civilization*. 11 Sep. 2012. <<http://cmcl.aai.uni-hamburg.de/>>.

“EpiDoc Guidelines.” *EpiDoc Guidelines*. 25 May 2013. <<http://www.stoa.org/epidoc/gl/dev/>>.

---. *EpiDoc: Epigraphic Documents in TEI XML*. 25 May 2013.
<<http://sourceforge.net/p/epidoc/wiki/Home/>>.

“Evertime: Antinoou.” *Evertime: Antinoou - A Standard Font for Coptic* 2012. 29 May 2013. <<http://www.evertime.com/fonts/coptic/>>.

Layton, Bentley. *A Coptic Grammar*. 3rd Edition, Revised. Wiesbaden: Harrassowitz, 2011. Print.

Orlandi, Tito. “The Library of the Monastery of Saint Shenute at Atripe.” *Perspectives on Panopolis: An Egyptian Town from Alexander the Great to the Arab Conquest*. Leiden: Brill, 2002. 211–231. Print.

Till, Walter C. “La séparation des mots en Copte.” *Bulletin de l’Institut français d’archéologie orientale* 60 (1960): 151–70.

Zeldes, Amir, Ritz, Julia, Lüdeling, Anke & Chiarcos, Christian “ANNIS: A Search Tool for Multi-Layer Annotated Corpora.” *Proceedings of Corpus Linguistics 2009* (2009) : n. pag. 10 Sep. 2012. <<http://ucrel.lancs.ac.uk/publications/cl2009/>>.

SCRIPTORIUM Part-of-Speech Tagsets for Sahidic Coptic

Amir Zeldes¹ & Caroline T. Schroeder²

1. Georgetown University

2. University of the Pacific

Version: 1.1.8_2016.08.31

1. Preamble

This document details guidelines for part-of-speech tagging Sahidic Coptic according to the SCRIPTORIUM project scheme. The tagging procedure assumes the text has already been normalized to the orthography and morpheme based segmentation described in the SCRIPTORIUM tokenization guidelines, which are closely related to the conventions found in Layton's (2004) grammar. In case of doubt we refer to Layton (2004) as well as Shisha-Halevy (1988).

As in all tagging projects, the aim is to achieve a practicable compromise between linguistic accuracy/usefulness, speed and reliability of human tagging, and performance of automatic tagging software. This means that in many cases concepts that are linguistically distinct are not distinguished since they are difficult to tell apart in practice in many cases, or determining some distinctions is too costly in terms of annotation time. Additionally, the project is using the CMCL lexicon, kindly provided by Prof. Tito Orlandi, which has its own, much more detailed scheme, so that in some cases the categories used here are chosen to be derivable from the CMCL scheme (see <http://cmcl.let.uniroma1.it/>).

There are two proposed tagsets, a coarse tagset with fewer tags for projects wishing to save annotation time, and a finer tagset with more detailed subcategories for some of the coarse grained tags, which is also expected to yield lower accuracy in automatic tagging. Links to the latest training models are provided from the SCRIPTORIUM website and have been tested and developed using the freely available TreeTagger (Schmid 1994, see <http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>).

2. Tagsets

The two tagsets described below are compatible with each other in that the fine-grained tagset uses the same overarching categories of the coarse one, but with further categories distinguished. The tag names are built 'hierarchically', so that additional letters in the name of a tag specify a special type of the superordinate category, e.g. all pronoun tags being with P, though not all tags with P are pronouns, as in PREP for prepositions.

In the coarse-grained list below, tags that have multiple fine-grained variants are followed by [*] (this is **not** part of the tag within the course-grained tagset).

Additionally, both tagsets admit certain cases where a single form contains two categories and must therefore be assigned two tags. This results in special underscore separated **portmanteau tags**, which are described in Section 2.3.

2.1 Coarse-Grained Tagset

Tag	Name	Examples
A[*]	Auxiliary tripartite base	α[ϰ], με[ϰ], τρε[ϰ], ...
ADV	Adverb	εβολ, ον, πως
ART	Article	π(ε), τ(ε), ν(ε), ζεν, κε
C[*]	Converter	ε, ετε, νε, ...
CONJ	Conjunction	αυω, η, μη, και, ειτε, ...
COP	Copula	πε/τε/νε
EXIST	Existential/possessive	ουν/μη
FM	Foreign material	παρα τουτο
FUT	Future	να
IMOD	Inflected modifier	τηρ[ϰ], ζωω[τ], ...
N[*]	Noun	αθητ, ρωμε, αρχη, ...
NEG	Negation	ν, αν, τη[σωτη]
NUM	Numeral	ογα, αναγ, ...
PDEM	Pronoun, demonstrative	πει/παι, τει/ται, νει/ναι
PINT	Pronoun, interrogative	ογ, νιη
PPER[*]	Pronoun, personal	ϰ,Ϸ,ι,†,ν,ανοκ,ανῖ,...
PPOS	Pronoun, possessive	πεϰ,τετῖ,πογ,πα,πωι,...
PREP	Preposition	ετβε, ζῖ, ν, ῖμο[ϰ], ...
PTC	Particle	δε, ῖσι, δε, ...
PUNCT	Punctuation	, , ' ...
UNKNOWN	Unknown morph, lacuna	в_ _ _ , _ _ _ _ , ...
V[*]	Verb	σωτη, σωτη, σωτη, ειρε, ο, αρι, ...
VBD	Verboid	νανου[ϰ], πεχα[ϰ], πεχε,...

2.2 Fine-Grained Tagset

For descriptions of the added fine-grained tags, marked in cursive type, see the coarse tag descriptions below.

<i>AAOR</i>	<i>AJUS</i>	<i>ANY</i>
<i>ACAUS</i>	<i>ALIM</i>	<i>AOPT</i>
<i>ACOND</i>	<i>ANEGAOR</i>	<i>APREC</i>
<i>ACONJ</i>	<i>ANEGJUS</i>	<i>APST</i>
<i>ADV</i>	<i>ANEGOPT</i>	<i>ART</i>
<i>AFUTCONJ</i>	<i>ANEGPST</i>	<i>CCIRC</i>

<i>CFOC</i>	NEG	PREP
<i>CPRET</i>	<i>NPROP</i>	PTC
CONJ	NUM	PUNCT
COP	PDEM	UNKNOWN
<i>CREL</i>	PINT	V
EXIST	<i>PPERI</i>	VBD
FUT	<i>PPERO</i>	<i>VIMP</i>
IMOD	<i>PPERS</i>	<i>VSTAT</i>
N	PPOS	

2.3 Portmanteau tags

In certain cases, one indivisible form corresponds to what normally constitutes two categories. This can happen either because of a phonological merger of two units, or because the formal marker of one category can be ‘zero’, i.e. have no form at all (usually in the case of 2nd person singular feminine forms). Portmanteau tags currently supported by the SCRIPTORIUM tools are:

tag	example	notes
AOPT_PPERS	ερε(σωτη)	Personal pron. within optative ε_ε. Note that ερε(σωτη) for 2nd pers. sg. fem. is also AOPT_PPERS, but nominal ερε(πρωμε σωτη) is only AOPT.
ACOND_PPERS	ερωαν(σωτη)	Personal pron. within conditional ε_ωαν. Note that ερωαν(σωτη) for 2nd pers. sg. fem. is also ACOND_PPERS, but nominal ερωαν(πρωμε σωτη) is only ACOND.
ACONJ_PPERS	τα(σωτη)	Truncated conjunctive 1st person (instead of ντασωτη)
ANEGPST_PPERS	μπε(σωτη)	Fused negative past 2nd pers. sg. fem. form.
APST_PPERS	αρ(σωτη)	Fused positive past 2nd pers. sg. fem. form.
CCIRC_PPERS	ερε(σωτη)	Fused circumstantial 2nd pers. sg. fem. form.
CFOC_PPERS	ερε(σωτη)	Fused focalized 2nd pers. sg. fem. form.
CPRET_PPERS	νερε(σωτη)	Fused preterit 2nd pers. sg. fem. form.
CREL_PPERS	ετρε(σωτη)	Fused relative 2nd pers. sg. fem. form.
IMOD_PPERO	μινιμο	The 2nd pers. sg. fem. form of ‘yourself’ (not to be confused with μινιμο(q) etc.)
PREP_PPERO	ερο	Any preposition where a 2nd pers. sg. fem. is realized as zero (not to be confused with ερο(q) etc.)
V_PPERO	(q)ντ	Verb forms with a fused 1st pers. sg. object, e.g. ντ ‘bring me’ from εινε ‘bring’, where the presuffixal form ντε is merged with the 1st pers. object marker -τ)

Note that in all cases, coarse grained tags can be substituted for fine grained ones, e.g. CCIRC_PPERS and CFOC_PPERS both become C_PPERS. Further combination tags are not ruled out and new ones will therefore be added if they are determined to be necessary.

2.4 Part of speech in conversion

In rare cases, a part of speech may appear in a syntactically unusual position. For example, an adverb or preposition may follow an article if they begin a phrase that is treated as a nominal phrase syntactically: the word $\epsilon\upsilon\lambda$ is tagged as an ADV, although in the sequence $\omicron\gamma|\epsilon\upsilon\lambda \ \xi\mu|\pi|\zeta\omega\mu\alpha$ ‘one (which is) out of the body’, it appears to behave like a noun. We consider such cases of ‘conversion’ between categories to be a syntactic phenomenon, and we therefore continue to tag $\epsilon\upsilon\lambda$ morphologically as an adverb.

An exception to this rule is the tagging of verbal infinitives following an article. In essence, almost any Coptic infinitive may be used as a noun, for example $\pi|\tau\omega\xi\mu$ ‘the call’. Cases such as these are widespread and are tagged as nouns, not as verbs, when the infinitive is used in this way.

3. Guidelines

The following guidelines describe the recommended assignment of part of speech tags to segmented morphemes. Fine-grained tags are given in the section describing the corresponding coarse-grained tag. In each example, the area corresponding to the tag under discussion is underlined. Vertical lines (‘pipes’) are used to segment morphemes for added clarity only.

3.1 Auxiliaries (A)

Auxiliaries include all conjugation bases in the tripartite patterns described in Layton (2004:251-290). These include both negative and positive variants and cover all lexical material preceding the subject noun or pronoun, e.g.:

- (1) $\underline{\alpha}|\eta|\zeta\omega\tau\bar{\mu}$ (3rd person masculine past tense)
- (2) $\underline{\alpha}\rho\epsilon|\zeta\omega\tau\bar{\mu}$ (2nd person feminine past tense, with zero subject)
- (3) $\underline{\mu}\pi|\eta|\zeta\omega\tau\bar{\mu}$ (negative past tense)

Note that when used with pronominal subjects, the optative and conditional conjugation encompass the subject pronoun, leading to a portmanteau tag like AOPT_PPERS (or A_PPERS in the coarse grained tagset):

- (4) $\underline{\epsilon}\eta\epsilon$ /AOPT_PPERS $\zeta\omega\tau\bar{\mu}$ (optative and 3rd pers. masc. pronoun)
- (5) $\underline{\epsilon}\eta\varpi\alpha\eta$ /ACOND_PPERS $\zeta\omega\tau\bar{\mu}$ (conditional and 3rd pers. masc. pronoun)

Fine-Grained Tags

The different individual fine-grained tags cover all distinct conjugation bases, making auxiliaries the largest fine-grained tag group. They are divided as follows:

APST	Auxiliary, past	α
ANEGPST	Auxiliary, negated past	$\bar{\mu}\pi(\epsilon)$

ANY	Auxiliary, ‘not yet’	ἤπατ(ε)
AAOR	Auxiliary, aorist	ᾠα, ᾠαρ(ε)
ANEGAOR	Auxiliary, negated aorist	με(ρε)
AOPT	Auxiliary, optative	ε[ϰ]ε, ερε
ANEGOPT	Auxiliary, negated optative	ῆνε
AJUS	Auxiliary, jussive	μαρ(ε)
ANEGJUS	Auxiliary, negated jussive	ἤπῑτρε
APREC	Auxiliary, precursive (‘after’)	ἤτερ(ε)
ACOND	Auxiliary, conditional	ε[ϰ]ᾠαν, ερᾠαν
ALIM	Auxiliary, limitative (‘until’)	ᾠαντ(ε)
ACONJ	Auxiliary, conjunctive	ῆ(τε)
AFUTCONJ	Auxiliary, future conjunctive	ταρ(ε)
ACAUS	Auxiliary, causative	τρε

Note that the irregular negation με in με-ᾠαε ‘it is not appropriate’ is also tagged NEG and not as ANEGAOR.

3.2 Adverbs (ADV)

Adverbs include indeclinable native Egyptian and Greek lexemes that modify verbs and other phrases as in the following examples.

- (6) τααῦζανε ἥμοϣ εματε/ADV ‘I shall glorify him greatly’
(7) πετ|ἥμαϣ/ADV ‘the one (who is) there’
(8) ἥπρμοϣ κακως/ADV ‘don’t die badly’

The first part of ‘complex prepositions’ is also tagged as an adverb, as in the following examples:

- (9) εβολ/ADV ῥῆ/PREP ‘from, out of’ (lit. ‘out in’)
(10) εροϣ/ADV ῥ/PREP ‘in towards’ (lit. ‘inside at’)

This does not apply to etymologically complex one-word prepositions derived e.g. from nouns for body parts (see the tag PREP for details), nor is the initial ε in words such as εβολ separated from the adverb (see segmentation guidelines).

3.3 Articles (ART)

Articles include definite articles, indefinite articles and article-like words such as κε/σε ‘other’. The following examples illustrate some variants:

- (11) π/ART ρωμε/N ‘the man’
(12) τε/ART κληρονομια/N ‘the inheritance’
(13) οϣ/ART νομος/N ‘a law’

- (14) ῥεν/ART ῥβηγε/N ‘(some) deeds’
 (15) κε/ART πονηρος/N ‘another wicked one’

Note that possessive pronouns like περ are not tagged as articles (see PPOS) and relative articles like πλετ are segmented to contain a relative converter (see C and CREL).

Articles followed by a noun beginning with Ϸ and consequently spelled θ or φ e.g. θε ‘the way’ are normalized and tokenized as τ and Ϸε before part-of-speech tagging, so that τ etc. can be tagged as an article alone (see segmentation guidelines).

3.4 Converters (C)

The class of converters, which is syntactically heterogeneous, is described in Layton (2004: 319-366). It includes four types of converters which have several realizations depending on their syntactic environment. In the coarse tagset, all converters are tagged as C, allowing for lower error rates in automatic tagging (especially by removing the distinction between circumstantial and relative conversions, which can be ambiguous). The examples below are for the four fine grained classes:

CCIRC	Converter, circumstantial	ε, ε[λ], ερε
CFOC	Converter, focalizing (a.k.a. 2 nd tenses)	ε, ερε, ετε, <u>ντ</u> [λ], <u>εντ</u> [λ]
CPRET	Converter, preterite	νε, νερε
CREL	Converter, relative	ετε, ετ, <u>ντ</u> [λ], <u>εντ</u> [λ], ετερε

Note that a following conjugation base is segmented separately from the converter (cf. segmentation guidelines), e.g.:

- (16) ντ/CREL λ/APST η|μογ ‘which he blessed’

The converter includes only ντ, while λ is a separate auxiliary base. The fused second person singular feminine form preceding a future marker, νερα (‘you(F) would’), is tokenized into norms and tagged as follows: νερ/CPRET_PPERS λ/FUT. Note that the normalized form of the future marker in this case remains λ, but the lemma is να.

3.5 Conjunctions (CONJ)

Conjunctions are indeclinable words of Greek and Egyptian origin which link phrases and clauses. No distinction is made between subordinating conjunctions which introduce clauses (‘because’, ‘lest’) and coordinating conjunctions which connect phrases (e.g. ‘and’, ‘or’).

- (17) αγω/CONJ αιειβεγ ‘and I became thirsty’
 (18) ειλω ἦμος ξε/CONJ μηποτε/CONJ ταειβε ‘saying [that:] lest I become thirsty’

In the first example, the coordinating conjunction $\alpha\gamma\omega$ ‘and’ appears. Note that it is still tagged as a conjunction even if the first coordinated phrase is missing. In the second example, two consecutive conjunctions appear: $\chi\epsilon$ ‘that, saying’ introduces the direct speech and the Greek origin $\mu\eta\pi\omicron\tau\epsilon$ ‘lest’ is a conjunction within the direct speech clause. Also note that the word $\chi\epsilon$, originally derived from $\chi\omega$ ‘say’ is not considered a verb in this usage.

3.6 Copulas (COP)

Copulas are markers in so-called nominal sentences which express predications of the sort A is B. The copula forms are $\pi\epsilon/\tau\epsilon/\nu\epsilon$. The tag COP is given also to copulas following a verbal clause for focalizing emphasis (i.e. ‘it is the case that...’), as illustrated below.

- (19) $\omicron\gamma\varsigma\alpha\epsilon\iota\mu$ $\pi\epsilon$ /COP ‘he is a doctor’
 (20) $\nu\epsilon\varphi\tau\omega\beta\varsigma$ $\mu\pi\chi\omicron\epsilon\iota\varsigma$ $\pi\epsilon$ /COP ‘(it is that) he prayed to God’

In the latter example, it is less obvious that $\pi\epsilon$ is the copula, as its predicate is formally a clause and the form never changes its gender or number (i.e. as $\tau\epsilon/\nu\epsilon$; this is also referred to as ‘invariable $\pi\epsilon$ ’). Though the English translation cannot convey the presence of the copula adequately, these types of cases are still tagged as COP (see Layton 2004:223).

3.7 Existentials (EXIST)

Existentials include the unique lexemes $\omicron\gamma\bar{\nu}$ and $\mu\bar{\nu}$ in both pure existential and possessive forms, positive and negative, illustrated in the following examples.

- (21) $\omicron\gamma\bar{\nu}$ /EXIST $\omicron\gamma\alpha$ $\epsilon\varphi\epsilon\iota\mu\epsilon$ $\mu\mu\omicron\kappa$ ‘there is one who is like you’
 (22) $\mu\bar{\nu}$ /EXIST $\gamma\bar{\mu}\gamma\alpha\lambda$ $\epsilon\varphi\chi\omicron\varsigma\epsilon$ $\epsilon\pi\epsilon\varphi\chi\omicron\epsilon\iota\varsigma$ ‘there is no servant who is above his master’

The same tag is also used for the indefinite durative present and the fixed phrase $\omicron\gamma\bar{\nu}$ $\beta\omicron\mu$ ‘be able’ literally ‘there is power’.

- (23) $\omicron\gamma\bar{\nu}\tau\alpha$ /EXIST ν /PPERO $\bar{\mu}\mu\alpha\gamma$ /ADV $\mu\pi\epsilon\bar{\nu}\epsilon\iota\omega\tau$ $\alpha\beta\bar{\rho}\alpha\gamma\alpha\mu$
 ‘we have Abraham our father’, lit. ‘exists to us ... of Abraham...’
 (24) $\bar{\mu}\mu\bar{\mu}$ /EXIST $\beta\omicron\mu$ $\nu\tau\epsilon$ | $\tau\epsilon$ | $\gamma\bar{\rho}\alpha\phi\eta$ $\beta\omega\lambda$ $\epsilon\beta\omega\lambda$ ‘scripture cannot be broken’

Note that the possessor pronoun is segmented apart from $\omicron\gamma\bar{\nu}\tau\alpha$ and tagged as a pronoun, and the accompanying $\bar{\mu}\mu\alpha\gamma$ is an adverb.

3.8 Foreign Material (FM)

Foreign material includes text that is lexically and syntactically from a foreign language. It is distinct from loan words. Loan words are lexical entries that originate in another language (e.g., Greek, Latin) but are used in Coptic with Coptic syntax. Foreign material

consists of words, especially multiword expressions, with foreign syntax. The writer has momentarily switched languages rather than embedded a loan word into a Coptic construction

- (25) ΟΥ ΠΑΡΑ ΤΟΥΤΟ/FM ΝΟΥ ΕΒΟΛ ΔΗ ΖΗΠΙΩΜΑ ΤΕ ‘it is therefore not part of the body’

3.9 Future Marker (FUT)

The future marker να, derived from the verb ‘go’ is not considered an independent verb form when introducing a second verb and marking future tense. The following example illustrates the construction.

- (26) † να/FUT ΖΟΥΤΕΚ ‘I will kill you’

In rare cases, forms other than να can be considered for the future marker, e.g. α in:

- (27) ΝΕΡ/CPRET_PPERS α/FUT αω ‘you would despise’ (2nd pers. fem.)

Contractions of multiple ν are usually restored in the normalization, so that a diplomatic sequence like τετναρπμεεγε ‘you will think’ are usually normalized and only then tagged as follows:

- (28) τετν/PPERS να/FUT ρ/V

3.10 Inflected modifiers (IMOD)

Inflected modifiers are a somewhat heterogeneous class of suffixally inflecting non-verboids, including the quantifier τηρ ‘all of’, the focus particle ογαα(τ) ‘only’ and the reflexive μμινμμο ‘oneself’ (see Layton 2004: 118-123 and contrast the tag VBD). The suffix itself is tokenized apart and tagged as PPERO. These items are tokenized apart even within larger phrases, as in the second examples below.

- (29) ανοκ ζωω/IMOD τ/PPERO ‘I, as for me / me too’
 (30) ε π τηρ/IMOD ι ‘in all of it, at all, wholly’

If the suffix is a 2nd pers. sg. fem. realized as zero, a portmanteau tag is assigned:

- (31) μμινμμο/IMOD_PPERO ‘yourself (2nd pers. sg. fem.)’

3.11 Nouns (N)

The tag N is used for all nouns, common and proper, though the fine-grained tagset offers the specific tag NPROP for proper nouns.

- (32) πεν ειωτ/N ‘our father’
 (33) αντωνιος/NPROP ‘Antonius’

Note that verbal infinitives in the durative patterns and elsewhere, though technically and etymologically nominal in nature, are nevertheless tagged as verbs in order to facilitate the retrieval of verbal lexemes across constructions.

- (34) † πιστεύε/V ἐπινοῦτε ‘I trust in God’

3.12 Negations (NEG)

The tag NEG is used for independent negative items that are not part of an auxiliary base. The following lexemes are given the tag NEG: ν, ἀν, τῆ and μῆ (negative imperative marker). The first two can occur in the same sentence, in which case one NEG tag is used for each. The third negates infinitives and is tokenized separately from the verb and surrounding auxiliaries. The fourth is also a separate token and is not considered a verb form or part of the verb εἶπε (this also applies to its lemmatization as an independent item, see lemmatization guidelines)

- (35) ᾤ/NEG ᾔνακληρονομεῖ ἡμῶν ἀν/NEG ‘he will not inherit you’
 (36) ἐγὼ ἀν τῆ/NEG σὺ τῆ ‘if they do not listen’
 (37) μῆ/NEG μοι κακῶς ‘don’t die badly!’

Note that the irregular negation με in με-οὔε ‘it is not appropriate’ is also tagged NEG and not as ANEGAOR.

3.13 Numerals (NUM)

The tag NUM is given to numerals and numerical constituents of complex numerals, as well as suffixed numerals as in the last example below.

- (38) πέν/NUM ἄρτοις ‘five (loaves) of bread’
 (39) ἑξήκ/NUM τέτ/NUM ‘twenty-four’
 (40) δι/NUM πάλιν/NUM ‘two times, twice’

Note that the indefinite article ὅς ‘a, one’ preceding a noun is tagged as ART, not NUM. Letters being used as numbers are considered NUM (including an alpha preceding a noun for the quantity ‘one’)

3.14 Demonstrative pronouns (PDEM)

The demonstrative pronouns, both attributive to the noun and substituting for a noun are tagged as PDEM.

- (41) ὅς/PDEM οὕτως ‘in this way’
 (42) ὅς/PDEM οὕτως ‘this is the way’

3.15 Interrogative pronouns (PINT)

This tag is used for the interrogative pronouns *οὔ* ‘what’, *νίμ* ‘who’, *τῶν* ‘where’, *ἅψ* ‘which’, *οὕηρ* ‘how much’. This is also true when they are used in complex phrases, as in the examples below.

- (43) *εἵτβε*/PREP *οὔ*/PINT ‘what for, why?’
(44) *εἵ*/PREP *τῶν*/PINT ‘where to?’

Note that the item *νίμ* is tagged PINT even when used after a noun to mean ‘some, any’.

3.16 Personal pronouns (PPER[*])

Personal pronouns generally receive the tag PPER, with three subtypes in the fine-grained subset for subject pronouns (PPERS), object pronouns (PPERO) and independent pronouns (PPERI).

- (45) *ἄ* *ψ*/PPERS *σώτῃ* *εἶπο* *κ*/PPERO ‘he heard you’
(46) *εἵτβηήτ* *ς*/PPERO ‘for her’

Note that ‘object’ pronouns include objects of prepositions and all suffixed pronouns except the subject markers of verboids of the type [*νᾱνοῦ*]*ψ*, [*πᾱχα*]*ψ* etc., which are tagged as PPERS.

- (47) *πᾱχα* *ψ*/PPERS ‘he said’

The independent pronouns are reserved for emphatic uses and nominal sentences, including nominal sentence subject forms like *ἄνῃ* ‘I’ and the full forms of the type *ἄνοκ* ‘I’.

- (48) *ἄνοκ*/PPERI *ζῶω* *τ*/PPERO *ἄνῃ*/PPERI *πᾱψ* *ζῃζᾱλ*
‘I, as for me, I am his servant’

Also note that possessive pronouns like *πᾱψ* ‘his’ are not segmented and receive a separate tag, PPOS.

3.17 Possessive pronouns (PPOS)

Much like demonstratives, all possessive pronouns, both attributive and standing in for a noun are tagged as PPOS. The personal suffix at the end of the pronoun is not separated, rather the entire forms, including *πᾱψ* ‘his’, *πᾱ* ‘my’ and ‘the one that belongs to’, *ποῦ* ‘your (fem.)’, *πῶν* ‘mine’ etc. The following example illustrates these different types of possessives:

- (49) *τᾱ*/PPOS *πᾱ*/PPOS *con* *τῶν*/PPOS *τᾱ* ‘the one of my brother is mine’

This tag only applies to prefixal, article-like possessives. Suffix possessives, such as πατ q ‘his foot’ are not tagged PPOS, but rather PPERO.

3.18 Prepositions (PREP)

This tag is used for all prepositions in both independent, prenominal states and presuffixal forms (which are tokenized apart from following suffixes). Note that prepositions that are historically derived from unverbized phrases but are now unsegmentable are tagged as one preposition, but complex prepositions involving a separable adverb are given two tags, ADV and PREP (cf. the tag ADV). Additionally, the *nota relationis* and accusative marker ν/ῖμο is regarded as a preposition. The following examples illustrate these principles.

- | | | |
|------|------------------|---|
| (50) | ἐτβε/PREP ογ | ‘for what? why?’ |
| (51) | εβολ/ADV ῖν/PREP | ‘from, out of’ (lit. ‘out in’) |
| (52) | εχῖν/PREP | ‘upon, on account of’ (from ‘to head of’) |

Also note that 2nd pers. sg. fem. objects often lead to portmanteau tags, e.g.:

- | | | |
|------|----------------|---------------------------------------|
| (53) | ῖμο/PREP_PPERO | ‘you (2nd pers. sg. fem. accusative)’ |
|------|----------------|---------------------------------------|

If in doubt as to whether a lexicalized combination is considered a single preposition, please refer to the formatted CMCL lexicon supplied with the project’s tokenization module. This lexicon will be updated with future versions of the guidelines to accommodate dubious cases as they arise.

3.19 Particles (PTC)

The class of particles contains all indeclinable words that do not belong to one of the other classes, most notably and frequently the apposition marker νγι ‘that is...’ and a large number of, mostly Greek origin, sentence modifying particles that tend to appear in the second, Wackernagel position as they do in Greek as well (e.g. δε, γαρ).

3.20 Punctuation (PUNCT)

All punctuation marks, including periods at any height in the line, commas (including punctuation added in editions when annotating edited texts) or even question marks, colons etc. if they are used, are all given the uniform tag PUNCT. If decorations are tokenized (tildes, clusters of dots etc.), they may also be tagged as PUNCT, though refer to the tokenization guidelines for recommendations on normalizing text before tagging.

3.21 Unknown, damaged and lost items (UNKNOWN)

The tag UNKNOWN is given to fragmentary word forms damaged or missing beyond the ability to reach a reliable part-of-speech assignment. It is understood in the case of larger lacunae that the string used to encode the visible part of a word may in fact contain

several words. In cases where it is clear where word divisions occur, multiple tokens with corresponding UNKNOWN tags are given.

- (54) ε[...]/UNKNOWN ‘?’
 (55) ε[...]/UNKNOWN π[...]/UNKNOWN ‘?’

Generally UNKNOWN tags are given even if the range of possibility is limited, i.e. even if we are certain a damaged morpheme is either an article or a possessive pronoun, an uncertain case is usually tagged as UNKNOWN.

3.22 Verbs (V[*])

The coarse tag V is given to all lexical verb forms that are not conjugation bases, also not including verboids, which receive a separate tag even in the coarse tagset due to their distinct syntax (see the tag VBD). In the fine-grained tagset, normal verb forms (V) are distinguished from stative verb forms (VSTAT) and imperatives (VIMP) as shown in the examples below. Note that verbal infinitives in the durative present are still tagged as verbs, although they are historically nominalized in this position, whereas nominalized infinitives following an article are understood as nouns, as in the last example. Verbs are tagged as VIMP only when they appear in the specific imperative form.

- (56) α υ σωτη̄/V ερο κ ‘he heard you’
 (57) † οβε/VSTAT ‘I am thirsty’
 (58) α.ι/VIMP ε ‘say it!’
 (59) ε̄ π π σοϋν/N ᾱ π νοϋτε ‘in the knowledge of God, the knowing of God’

Also note that in rare cases, object pronouns that are realized as zero will lead to portmanteau tags, e.g.:

- (60) τετην/PPERS ντ/V_PPERO ‘you bring me’

Since ντ= as the presuffixal form of εινε ends in τ, the object pronoun -τ ‘me’ is subsequently dropped. However the portmanteau tag reflects the presence of a grammatical object.

For compound verbs (see §180 in Layton), the entire compound is considered “a single unit in boundness, syntax, and meaning.” Therefore, the entire compound is tagged V. The components of the compound may be annotated further on a morph level annotation. (See Transcription guidelines for more information on bound groups, morphemes, and word segmentation.) Common examples include compound verbs formed with †-, p-, and αι-.

- (61) ετ/CREL πνοβε/V
 (62) ε/CCIRC κ/PPERS †CBΩ/V

The basic criterion for identifying compound verbs is the absence of an article: $\rho\nu\omicron\upsilon\epsilon$ ‘to sin’ is considered as single, compound verb (which can still be analyzed morphologically into two units $\rho+\nu\omicron\upsilon\epsilon$, perhaps like English ‘sin-ify’, if there were such a word). However $\rho\ \pi\ \mu\epsilon\epsilon\gamma\epsilon$ ‘to think’ looks exactly like any verb + definite noun phrase combination, and is therefore tagged as three units despite being a common lexicalized combination: it comprises a verb, an article and a noun.

Exceptions: Some object nouns cannot appear as definite, or are made definite other than by an article. These include objects with $\nu\iota\mu$ ‘some, any’, $\lambda\alpha\lambda\gamma$ ‘something’ and $\rho\omicron\iota\mu\epsilon$ ‘some (ones)’, number words, as well as verbal objects with a suffixal possessive pronouns, such as $\kappa\epsilon\ \rho\alpha\tau\ \tau\ \upsilon$ ‘set one’s foot’ (the foot is definite). Even though they may appear next to a verb without an article, these are tokenized and tagged apart from the verb (for possessed objects, the possessive is its own token, tagged PPERO, not PPOS).

Compound verbs containing a specific imperative form are also considered VIMP:

- (63) $\alpha\rho\iota\epsilon\mu\epsilon\lambda\lambda$ /VIMP ‘serve!’ (imperative of compound $\rho\epsilon\mu\epsilon\lambda\lambda$)

3.23 Verboids (VBD)

The category VBD is given to a small class of suffixally inflected predicates described in Layton (2004: 297-304), including the common $\mu\epsilon\chi\epsilon$ -/ $\mu\epsilon\chi\alpha\epsilon$ ‘say’, $\mu\alpha\lambda\omicron\upsilon\gamma\epsilon$ ‘be good’ etc., but not including possessive existentials of the type $\omicron\gamma\iota\tau\epsilon$ - (see the tag EXIST). The personal suffix following a VBD is tagged as its subject, i.e. PPERS (or simply PPER in the coarse tagset).

- (64) $\mu\epsilon\chi\alpha$ /VBD τ /PPERS ‘he said’
 (65) $\mu\alpha\lambda\omicron\upsilon\gamma$ /VBD ς /PPERS ‘she/it is good’

For the form $\mu\epsilon\omega\alpha\epsilon$ note that two analyses exist. When it is declinable and literally means ‘X does not know’ (also prenominal $\mu\epsilon\omega\epsilon$ -), then it is VBD. When it is the lexicalized adverb form $\mu\epsilon\omega\alpha\kappa$ meaning ‘maybe’ (etymologically from ‘you never know’), it is a single unit, tagged ADV. Note that the latter form does not agree with the addressee if they are not masculine singular. Contrast the following examples from Layton (2004:303):

- (66) $\mu\epsilon\omega\alpha$ /VBD τ /PPERS $\mu\eta\mu\alpha\gamma\ \epsilon\tau\epsilon\mu\alpha\kappa\omega\ \mu\epsilon\omega\tau\ \mu\eta\kappa\omicron\varsigma\mu\omicron\varsigma$
 ‘he does not know when he will leave the world’
 (67) $\mu\epsilon\omega\alpha\kappa$ /ADV $\tau\ \mu\alpha\sigma\omega\ \chi\alpha\tau\epsilon\tau\eta\gamma\tau\iota$ ‘maybe I’ll stay with you’

In the latter example, the addressee is plural ($\tau\eta\gamma\tau\iota$), but the form remains $\mu\epsilon\omega\alpha\kappa$, indicating that it is an unanalyzed adverb.

4. References

- Layton, Bentley (2004), *A Coptic Grammar*. Second Edition, Revised and Expanded. (Porta linguarum orientalium 20.) Wiesbaden: Harrassowitz.
- Schmid, Helmut (1994), Probabilistic part-of-speech tagging using decision trees. *Proceedings of the Conference on New Methods in Language Processing*. Manchester, UK, 44–49.
- Shisha-Halevy, Ariel. 1988. *Coptic Grammatical Chrestomathy. A Course for Academic and Private Study*. (Orientalia Lovaniensia Analecta 30.) Leuven: Peeters.

Coptic SCRIPTORIUM – Lemmatization Guidelines

Version 1.1.0 / 2016-11-02

Amir Zeldes

Preamble

The purpose of lemmatization is to facilitate finding variant and inflected forms that are related to the same lexical entry, roughly equivalent to a dictionary entry. However in many cases, it may be unclear what the underlying, uninflected form of a word is: is the lemma of the pronoun ‘me’ defined as ‘I’ (i.e. the nominative form)? Should the lemma of ‘us’ then be ‘we’? Alternatively we could put all personal pronouns under one lemma: then ‘we’, ‘us’, ‘I’, and ‘me’ all belong to the same lemma, but which form should be taken for the common lemma?

There can be many arguments for and against certain practices. In these guidelines we attempt to give a set of instructions for Coptic which is: a. easy to apply consistently and b. useful for searching purposes.

Guidelines by Part-of-Speech Class

Articles and copulas

Articles are lemmatized according to the non-assimilated, simple short form of the corresponding masculine singular article (if distinct). This means that the lemma of π, πε, τ, τε, η, ηε and ι (assimilated form of η before a labial consonant) is for all of the above π. For indefinites ογ and γεν there is no special masculine form, but the singular lemma ογ is taken for the plural γεν and also for the variant spelling γ.

Copulas follow a similar rule: the lemma for all three number/gender forms (πε/τε/ηε) is πε.

Pronouns

Personal pronouns

Lemmas are mainly helpful where they deliver added value over searching for plain strings. It is therefore useful to give common lemmas for each of the personal forms: first person, second M/F, ... Given that the SCRIPTORIUM part-of-speech guidelines already distinguish subject and object pronouns, it is considerably more useful to group subjects and objects of the same person together, while not distinguishing the different forms (e.g. ⲧ, ⲓ for first person) which can be found using a plain-text search anyway. We therefore annotate the following personal pronouns (SCRIPTORIUM tags in PPER*, i.e. PPERs,

PPERO, PPERI) with the following lemmas, based on the independent stressed pronoun forms (note that lemmatization is based on normalized forms without supralinear strokes or other diacritics; cf. transcription and normalization guidelines):

Person	Lemma	Pronoun forms
1st sg.	ΔΝΟΚ	ΔΝΟΚ, ΔΝΓ, †, Ι, ΝΤ, Τ, Δ
2nd sg. masc.	ΝΤΟΚ	ΝΤΟΚ, ΝΤΚ, Κ, Γ, ΤΚ
2nd sg. fem.	ΝΤΟ	ΝΤΟ, ΝΤΕ, ΤΕ, ΤΡ, Ρ, Ε
3rd sg. masc.	ΝΤΟQ	ΝΤΟQ, Q
3rd sg. fem.	ΝΤΟC	ΝΤΟC, C
1st pl.	ΔΝΟΝ	ΔΝΟΝ, ΔΝ, Ν, ΤΝ, CΝ
2nd pl.	ΝΤΩΤΝ	ΝΤΩΤΝ, ΝΤΕΤΝ, ΤΝ, ΤΗΥΤΝ
3rd pl.	ΝΤΟΟΥ	ΝΤΟΟΥ, Υ, ΟΥ, CΕ, CΟΥ

The pronoun lemmas alone therefore primarily give access to search by person (1st, 2nd ...); to cross-reference these with the form, e.g. independent pronoun, cross-reference the POS annotation (in ANNIS: pos="PPERI"). For a specific subform (e.g. ΔΝΓ not ΔΝΟΚ) use the form search norm="ΔΝΓ".

Possessives, interrogatives and demonstratives

Interrogative pronouns are each equivalent to their own lemma, i.e. οΥ is lemmatized οΥ and ΝΙΜ as ΝΙΜ.

Possessive, and demonstrative pronouns are lemmatized to their own normalized form, but with one modification: non-masculine singular determiners are given the masculine form, i.e. the lemma of πεq is πεq, the lemma of πα is πα etc., but the lemma of τεc and nec is also nec. Similarly, the lemma of πει and νει is πει, and the lemma of παι and ται is παι. This allows an easier search for all possessives (in ANNIS: pos="PPOS", finds πεγ, τεγ, νογ ...), all third person plural possessives (lemma="πεγ", finds πεγ, τεγ and νεγ) and all third person plural possessives of feminine objects (norm="νεγ"), and similarly for demonstratives.

Adverbs, particles and conjunctions

Adverbs, particles and conjunctions are always given their own normalized form as a lemma. This includes Greek adverbs in -ωc, which are lemmatized as such, e.g. ζωλωc has the lemma ζωλωc.

Nouns

Nouns are given their dictionary form as a lemma. For most nouns, singular and plural forms are identical, meaning there is no dilemma. For nouns with irregular plural forms, the singular form is taken as a lemma, e.g. ζωβ 'deed' is the lemma of both singular ζωβ and plural γβηγε, and similarly, possessed forms like τοοτ(q) are lemmatized to the absolute form, i.e. τωπε. In order to find irregular forms, one can then simply search for

nouns whose lemma is different from the noun form (in ANNIS: lemma != norm). The same rules apply to proper nouns, though these rarely occur in the plural.

For nouns which only occur in the possessed form, if both prenominal and presuffixal forms exist, the prenominal is taken as the lemma, e.g. ζῆλᾱς and ζῆνε- ‘(one’s) will’ are lemmatized as ζῆνε. If only a presuffixal form exists, it is taken as the lemma as well, e.g. ἡλιατῶ ‘blessed is...’ has the lemma ἡλιατ.

Nouns that have related masculine and feminine forms are considered separate lemmas. For instance, the noun υἱος ‘son’ is its own lemma, and the separate noun θυγατήρ ‘daughter’ also has a separate lemma (which is θυγατήρ). Similarly, Greek words in -ος are considered separate from related words in -ον, e.g. πονηρός ‘wicked person’ is its own lemma, and so is the separate πονηρόν ‘wicked deed/thing’ an independent lemma.

Verbs

Verbs are lemmatized to the form of the absolute infinitive. This means that special prenominal or presuffixal forms are lemmatized to their respective dictionary entries, e.g. ὀφείλω and ὀφείλ- are lemmatized as ὀφείλω ‘choose’. The same applies to stative and imperative forms, which are lemmatized to the dictionary entry, e.g. κητῆς has κωτ as a lemma and ἀπὶ has εἶπε. Likewise for prenominal forms, ἔστι and εἶ is lemmatized as ὄν and εἶπε.

Note that auxiliaries are not lemmatized to their etymological verbs, i.e. the lemma of the past tense ἔσθι is not εἶπε but ἔσθι. Additionally, the negative imperative marker μή is lemmatized as μή as well, as it is considered to be a form of negation independent from the verb εἶπε. However, the negative imperative of εἶπε itself, μὴ εἶπε IS lemmatized as εἶπε (since it is a morphological imperative of εἶπε itself, and functions as part of its paradigm with the sense ‘to do’).

For fused verb-object forms like ἔλθω ‘bring me’, see Portmanteau Tags.

Prepositions

Prepositions are lemmatized to their standard form **before noun phrases**. Therefore the lemma of ἐν- and ἐπὶ- is ἐν. For preposition forms containing a second person singular feminine pronoun (realized as zero), e.g. ἐν σοὶ ‘on you (fem.)’, ὀπίσθω ‘behind you (fem.)’ etc. see Portmanteau Tags.

Existential and possessive predicates

The existential predicates are lemmatized as οὐκ ‘there is’ and οὐκ ‘there isn’t’ (again note that lemmatization does not contain supralinear strokes). Like auxiliaries, the related possessive predicates are lemmatized using their form before a nominal subject: οὐκ ὄντι

and **мнтѣ**. Note that forms with two **м**, such as **ммн** and **ммнтѣ**- etc. are considered to be orthographic variants of the forms with one **м**, and should therefore be not only lemmatized with one **м**, but also normalized to such forms.

Auxiliaries, negations and future marker

Auxiliaries are generally lemmatized to their form when preceding a nominal subject. Attention should be paid to auxiliaries sometimes ending in **-ѣ**: in normalized orthography, this is generally present before a nominal subject. The lemma of **маꙗѣ**- and **маꙗ**- (jussive) is **маꙗѣ**, and the lemma of **ѡант**- and **ѡантѣ** is **ѡантѣ**.

However, the lemmas of auxiliaries that sometimes contain an intermediate pronoun do not contain that pronoun when they occur uninterrupted: the lemmas of **ѣꙗꙗн** (conditional) and **ѣꙗѣ** (optative) remain **ѣꙗꙗн** and **ѣꙗѣ**. These receive the tags **ACOND** and **AOPT** respectively. For cases with an intervening pronoun, which receive different tags, see Portmanteau Tags.

Negative morphemes such as **н**, **ан** and **тн** are their own lemmas (the form **м** before a labial is also lemmatized as **н**). The negative imperative marker **мнꙗ** is lemmatized as itself (**мнꙗ**), and NOT as **ѣꙗѣ**.

The future marker is given its own lemma **на**. Note that the lemma remains so whenever a future marker is separately identified, even if the diplomatic realization is assimilated and reduced to **а**, e.g. in complex forms like **тетна** ‘you will... (pl.)’ or **неꙗ** ‘you would have (fem. sg.)’.

Converters

Like auxiliaries, converters are lemmatized to their form before a nominal subject, viz.:

CCIRC/CFOC:	ѣꙗѣ
CREL:	ѣтѣꙗѣ
CPRET:	неꙗѣ

For second person singular feminine **ѣꙗ/ѣꙗѣ** (lemma="**ѣꙗ_нто**") see Portmanteau Tags.

Inflected modifiers

Modifiers of the type **ꙗꙗꙗ-**, **ммнммм-**, **маꙗꙗꙗ-**, **тнꙗ-** are lemmatized to their form before the **third person masculine singular** pronoun **ѣ**. Thus **ммнммм-** and **ммнмммꙗ-** are lemmatized as **ммнммм**. The portmanteau form **ммнммм** (yourself, fem. sg.) is lemmatized **ммнммм_нто** (see Portmanteau Tags).

Portmanteau Tags

Some fused items receive a so-called portmanteau tag representing two categories at once. For example, the form $\epsilon\varphi\omega\alpha\lambda\iota$ is considered to contain a conditional auxiliary and a subject pronoun: $\text{pos}=\text{"ACOND_PPERS"}$. In order to facilitate finding such cases regardless of the pronoun in use, in tags containing a conjugation base and a personal pronoun the form is lemmatized using both lemmas, separated by an underscore. For example, the lemmas of $\epsilon\iota\omega\alpha\lambda\iota$, $\epsilon\varsigma\omega\alpha\lambda\iota$ and $\epsilon\varphi\omega\alpha\lambda\iota$ are $\epsilon\rho\omega\alpha\lambda\iota_a\lambda o\kappa$, $\epsilon\rho\omega\alpha\lambda\iota_n\tau o\varsigma$ and $\epsilon\rho\omega\alpha\lambda\iota_n\tau o\iota$ respectively. The lemma of $\epsilon\rho\omega\alpha\lambda\iota$ remains $\epsilon\rho\omega\alpha\lambda\iota$ ($\text{pos}=\text{"ACOND"}$), unless it contains a second person feminine singular subject, in which case the lemma is $\epsilon\rho\omega\alpha\lambda\iota_n\tau o$ according to the rule above.

For the past tense second person singular feminine form $\alpha\rho$ the lemma is similarly $\alpha_n\tau o$ ($\text{pos}=\text{"APST_PPERS"}$). The form $\mu\mu\iota\mu\iota\mu\iota\mu\iota$ (yourself, fem. sg.) is identical to the base of other personal forms, but is lemmatized $\mu\mu\iota\mu\iota\mu\iota\mu\iota_n\tau o$ just like other forms containing a personal pronoun.

The same principle applies to prepositions: forms containing a second person singular feminine pronoun (realized as zero) are given portmanteau lemmas, e.g. $\epsilon\chi\omega$ 'on you (fem.)' has $\epsilon\chi\iota_n\tau o$, $\kappa\omega$ 'behind you (fem.)' has $\kappa\alpha_n\tau o$ etc.

For circumstantial or focalizing converter + second person feminine singular, the lemma $\epsilon\rho\epsilon_n\tau o$ is used (and similarly preterit $\kappa\epsilon\rho\epsilon_n\tau o$ and relative $\epsilon\tau\epsilon\rho\epsilon_n\tau o$).

Verbs containing an object pronoun, such as $\eta\tau$ 'bring me' are lemmatized using the base form of the verb and the pronoun's lemma: $\epsilon\iota\kappa\epsilon_a\lambda o\kappa$.

FAQ

What is the lemma of...

Form	Norm	Lemma	Notes
$\chi\omega\omega(\iota)$	$\chi\omega\omega$	$\chi\omega\omega$	
$\chi\epsilon\kappa\alpha\varsigma$	$\chi\epsilon\kappa\alpha\alpha\varsigma$	$\chi\epsilon\kappa\alpha\alpha\varsigma$	just a normalization issue
$\kappa\alpha\tau\alpha\rho o(\iota)$	$\kappa\alpha\tau\alpha\rho o$	$\kappa\alpha\tau\alpha$	
$\mu\mu\iota\mu\iota\tau\alpha(\iota)$	$\mu\mu\iota\mu\iota\tau\alpha$	$\mu\mu\iota\mu\iota\tau\epsilon$	
$\omicron\gamma\epsilon\iota$	$\omicron\gamma\epsilon\iota$	$\omicron\gamma\alpha$	'one' (fem.) is considered an inflected form
$\sigma\iota\nu o\gamma\omicron o\mu$	$\sigma\iota\nu o\gamma\omicron o\mu$	$\sigma\iota\nu o\gamma\omega\mu$	complex plural still gets lemmatized as sg.
$\varsigma\alpha\beta o$	$\varsigma\alpha\beta o$	$\varsigma\beta o$	
$\tau\varsigma\alpha\beta o$	$\tau\varsigma\alpha\beta o$	$\tau\varsigma\alpha\beta o$	note this is one norm unit
$\varsigma\gamma\eta\kappa\epsilon\rho\alpha$	$\varsigma\gamma\eta\kappa\epsilon\rho\alpha$	$\varsigma\gamma\eta\kappa\epsilon\rho\alpha$	Greek words should normalize η for η

(π)τηρϚ	τηρϚ	τηρϚ	when used as a noun, this is one norm
ναιδτ(Ϛ)	ναιδτ	ναιδτ	
νι	νι	πι	has a separate lemma from πει/τει/νει

What about units with complex morphology?

The morph layer is not a type of lemmatization. We lemmatize the whole word, which may contain inflected forms of stems below the word level. Thus the lemma of ρχρεια is ρχρεια, not ειρερεια.