

# SCRIPTORIUM Part-of-Speech Tagsets for Sahidic Coptic

Amir Zeldes<sup>1</sup> & Caroline T. Schroeder<sup>2</sup>

1. Humboldt-Universität zu Berlin

2. University of the Pacific

*Version:* 1.0.1\_2013.7.6a

## 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 coarse-grained tagset).

## 2.1 Coarse-Grained Tagset

Tag	Name	Examples
A[*]	Auxiliary tripartite base	α[ϰ], με[ϰ], τρε[ϰ], ...
ADV	Adverb	εβολ, οη, πωσ
ART	Article	π(ε), τ(ε), ν(ε), ρεν, κε
C[*]	Converter	ε, ετε, νε, ...
CONJ	Conjunction	αγω, η, μη, και, ειτε, ...
COP	Copula	πε/τε/νε
EXIST	Existential/possessive	ογν/μν
FUT	Future	να
IMOD	Inflected modifier	τηρ[ϰ], ρωω[τ], ...
N[*]	Noun	αθητ, ρωμε, αρχη, ...
NEG	Negation	ν, αν, τμ[σωτμ]
NUM	Numeral	ογα, συναγ, ...
PDEM	Pronoun, demonstrative	πει/παι, τει/ται, νει/ναι
PINT	Pronoun, interrogative	ογ, νιμ
PPER[*]	Pronoun, personal	ϰ, ϰ, ι, †, ν, ανοκ, αν̄, ...
PPOS[*]	Pronoun, possessive	πεϰ, τετ̄ν, πογ, πα, πωι, ...
PREP	Preposition	ετβε, ρ̄ν, η, ἡμο[ϰ], ...
PTC	Particle	δε, ν̄σι, δε, ...
PUNCT	Punctuation	. , ' ...
UNKNOWN	Unknown morph, lacuna	β_ _ _ , _ _ ος, _ _ _ , ...
V[*]	Verb	σωτμ, σωτπ, σοτπ, ειρε, ο, αρι, ...
VBD	Verboid	νανογ[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>ANY</i>	<i>FUT</i>
<i>ACAUS</i>	<i>AOPT</i>	<i>IMOD</i>
<i>ACOND</i>	<i>APREC</i>	<i>N</i>
<i>ACONJ</i>	<i>APST</i>	<i>NEG</i>
<i>ADV</i>	<i>ART</i>	<i>NPROP</i>
<i>AFUTCONJ</i>	<i>CCIRC</i>	<i>NUM</i>
<i>AJUS</i>	<i>CFOC</i>	<i>PDEM</i>
<i>ALIM</i>	<i>CPRET</i>	<i>PINT</i>
<i>ANEGAOR</i>	<i>CONJ</i>	<i>PPERI</i>
<i>ANEGJUS</i>	<i>COP</i>	<i>PPERO</i>
<i>ANEGOPT</i>	<i>CREL</i>	<i>PPERS</i>
<i>ANEGPST</i>	<i>EXIST</i>	<i>PPOS</i>

PREP	UNKNOWN	VIMP
PTC	V	VSTAT
PUNCT	VBD	

### 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) α|q|εωτῃ (3rd person masculine past tense)
- (2) αρε|εωτῃ (2nd person feminine past tense, with zero subject)
- (3) ῃπ|ι|εωτῃ (negative past tense)

Note that when used with pronominal subjects, the optative and conditional conjugation bases are split around the subject pronoun, leading e.g. to two separate epsilons receiving the auxiliary tag (one tag for each segmented epsilon), and likewise for the conditional morphemes:

- (4) ε|q|εωτῃ (3rd person masculine optative, two auxiliary tags)
- (5) ε|q|ωαν|εωτῃ (3rd person masculine conditional, two auxiliary tags)

#### *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	ῃπ(ε)
ANY	Auxiliary, ‘not yet’	ῃπατ(ε)
AAOR	Auxiliary, aorist	ωα, ωαρε
ANEGAOR	Auxiliary, negated aorist	με(ρε)
AOPT	Auxiliary, optative	ε[q]ε, ερε
ANEGOPT	Auxiliary, negated optative	ῃνε
AJUS	Auxiliary, jussive	μαρ(ε)
ANEGJUS	Auxiliary, negated jussive	ῃπῑτρε
APREC	Auxiliary, precursive (‘after’)	ῃτερ(ε)

ACOND	Auxiliary, conditional	ε[υ]ϰαν, ερϰαν
ALIM	Auxiliary, limitative ('until')	ϰαντ(ε)
ACONJ	Auxiliary, conjunctive	ν̄(τε)
AFUTCONJ	Auxiliary, future conjunctive	ταρ(ε)
ACAUS	Auxiliary, causative	τρε

### 3.2 Adverbs (ADV)

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

- (6) τααϣζανε ἡμοϣ εματε/ADV 'I shall glorify him greatly'  
(7) πετ|ἡμαϣ/ADV 'the one (who is) there'  
(8) ἡπρμοϣ 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 ϣν̄/PREP 'from, out of' (lit. 'out in')  
(10) εροϣν/ADV ϣι/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 ϣωμε/N 'the man'  
(12) τε/ART κληρονομια/N 'the inheritance'  
(13) οϣ/ART νομος/N 'a law'  
(14) ϣεν/ART ϣβηγε/N '(some) deeds'  
(15) κε/ART πονηρος/N 'another wicked one'

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

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



- (19) ογcαειν πe/COP ‘he is a doctor’  
 (20) νεqτωβq μπcοειc πe/COP ‘(it is that) he prayed to God’

In the latter example, it is less obvious that πe is the copula, as its predicate is formally a clause and the form never changes its gender or number (i.e. as τε/νε; this is also referred to as ‘invariable πe’). 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 ογñ and μñ in both pure existential and possessive forms, positive and negative, illustrated in the following examples.

- (21) ογñ/EXIST ογα εφεινε μμοκ ‘there is one who is like you’  
 (22) μñ/EXIST q̃μq̃αλ εq̃οce επεq̃οειc ‘there is no servant who is above his master’

The same tag is also used for the indefinite durative present and the fixed phrase ογñ σομ ‘be able’ literally ‘there is power’.

- (23) ογñτα/EXIST n/PPERO μμαγ/ADV μπeneιδωτ αβραq̃αμ  
 ‘we have Abraham our father’, lit. ‘exists to us ... of Abraham...’  
 (24) μμñ/EXIST σομ nτε|τε|γραφη βωλ εβολ ‘scripture cannot be broken’

Note that the possessor pronoun is segmented apart from ογñτα and tagged as a pronoun, and the accompanying μμαγ is an adverb.

### 3.8 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.

- (25) † να/FUT q̃οτβεκ ‘I will kill you’

### 3.9 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 μμινμμō= ‘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.

- (26) ανοκ q̃ωω/IMOD τ/PPERO ‘I, as for me / me too’  
 (27) ε π τηρ/IMOD q ‘in all of it, at all, wholly’

### 3.10 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.

- (28) ΠΕΝ ΕΙΩΤ/N 'our father'  
(29) 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.

- (30) † ΠΙΣΤΕΥΕ/V ΕΠΙΝΟΥΤΕ 'I trust in God'

### 3.11 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, ΑΝ, Τῃ. The former two tags can occur in the same sentence, in which case one NEG tag is used for each. The latter tag negates infinitives and is tokenized separately from the verb and surrounding auxiliaries.

- (31) ἢ/NEG ἤΝΑΚΛΗΡΟΝΟΜΕΙ ἡΜΟΚ ΔΗ/NEG 'he will not inherit you'  
(32) ΕΥΦΑΝ Τῃ/NEG ΣΩΤῃ 'if they do not listen'

### 3.12 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.

- (33) ΠΟΥ/NUM ἥοεικ 'five (loaves) of bread'  
(34) ΧΟΥΤ/NUM ΔΥΤΕ/NUM 'twenty-four'  
(35) ἡ|ΣΕΠ ΣΝΔΥ/NUM 'two times, twice'

Note that the indefinite article οὐ 'a, one' preceding a noun is tagged as ART, not NUM.

### 3.13 Demonstrative pronouns (PDEM)

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

- (36) Ν ΤΕΙ/PDEM ζε 'in this way'  
(37) ΤΑΙ/PDEM ΤΕ Τ ΖΕ 'this is the way'

### 3.14 Interrogative pronouns (PINT)

This tag is used for the interrogative pronouns οὐ 'what', ΝΗ 'who', ΤΩΝ 'where', ΔΩ 'which', ΟΥΗΡ 'how much'. This is also true when they are used in complex phrases, as in the examples below.

- (38) ετβε/PREP ου/PINT ‘what for, why?’  
 (39) ε/PREP των/PINT ‘where to?’

### 3.15 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).

- (40) α υ/PPERS σωτῆ ερο κ/PPERO ‘he heard you’  
 (41) ετβηητ ε/PPERO ‘for her’

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

- (42) πεχα υ/PPERS ‘he said’

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

- (43) ανοκ/PPERI ζωω τ/PPERO ανῑ/PPERI πεϣ ζῆχαλ  
 ‘I, as for me, I am his servant’

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

### 3.16 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 πεϣ ‘his’, πα ‘my’ and ‘the one that belongs to’, πογ ‘your (fem.)’, πωι ‘mine’ etc. The following example illustrates these different types of possessives:

- (44) ταν/PPOS παν/PPOS συν τωι/PPOS τε ‘the one of my brother is mine’

### 3.17 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 univerbized phrases but are now unsegmentable are tagged as one preposition, but complex preposition involving a separable adverb are given two tags, ADV and PREP (cf. the tag ADV). Additionally, the



*nota relationis* and accusative marker  $\kappa/\bar{\kappa}\mu\sigma$  is regarded as a preposition. The following examples illustrate these principles.

- |   |   |
|---|---|
| (45) $\epsilon\tau\beta\epsilon$ /PREP $\omicron\gamma$           | ‘for what? why?’                          |
| (46) $\epsilon\beta\omicron\lambda$ /ADV $\gamma\bar{\eta}$ /PREP | ‘from, out of’ (lit. ‘out in’)            |
| (47) $\epsilon\chi\bar{\eta}$ /PREP                               | ‘upon, on account of’ (from ‘to head of’) |

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.

### 3.18 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  $\kappa\omicron\iota$  ‘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.  $\Delta\epsilon$ ,  $\gamma\alpha\rho$ ).

### 3.19 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 in the context of normalized text.

### 3.20 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.

- |   |     |
|---|-----|
| (48) $\epsilon[...]$ /UNKNOWN                     | ‘?’ |
| (49) $\epsilon[...]$ /UNKNOWN $\pi[...]$ /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.

### 3.21 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.

- (50) ⲁ ϣ ⲩⲱⲧⲙ̄/V ⲉⲣⲟ ⲕ      ‘he heard you’  
 (51) ⲛ ⲟⲃⲉ/VSTAT      ‘I am thirsty’  
 (52) ⲁⲭⲓ/VIMP ⲥ      ‘say it!’  
 (53) ⲉ̄ⲙ ⲡ ⲩⲟⲟϥⲛ/N ⲙ̄ ⲡ ⲛⲟϥⲧⲉ      ‘in the knowledge of God, the knowing of God’

### 3.22 Verboids (VBD)

The category VBD is given to a small class of suffixally inflected predicates described in Layton (2004: 297-304), including the common ⲡⲉⲭⲉ-/ⲡⲉⲭⲁⲣ ‘say’, ⲛⲁⲛⲟϥⲣ ‘be good’ etc., but not including possessive existentials of the type ⲟϥⲛⲧⲉ- (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).

- (54) ⲡⲉⲭⲁ/VBD ϣ/PPERS      ‘he said’  
 (55) ⲛⲁⲛⲟϥ/VBD ⲥ/PPERS      ‘she/it is good’

## 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.