

# SCRIPTORIUM Part-of-Speech Tagsets for Sahidic Coptic

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

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[σωτμ]
NUM	Numeral	ογα, συναυ, ...
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	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>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 2<sup>nd</sup> 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 εβολ is tagged as an ADV, although in the sequence ογ|εβολ χμ|π|σωμα ‘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 εβολ 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 π|τωχμ ‘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) α|q|σωτῃ (3rd person masculine past tense)
- (2) αpε|σωτῃ (2nd person feminine past tense, with zero subject)
- (3) ῃπ|ι|σωτῃ (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) εqε/AOPT\_PPERS σωτῃ (optative and 3rd pers. masc. pronoun)
- (5) εqαν/ACOND\_PPERS σωτῃ (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	ῃπ(ε)

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	ἠπῖτρε
APREC	Auxiliary, precursive (‘after’)	ἦτερ(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) εροϣν/ADV ῥν/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 <sup>nd</sup> 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 q|cmoy ‘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 αλειβεγ ‘and I became thirsty’  
 (18) ειλω ἡμος 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  $\nu$ /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  $\beta\omicron\mu$   $\nu\tau\epsilon|\tau\epsilon|\gamma\bar{\rho}\alpha\phi\eta$   $\bar{\nu}\omega\lambda$   $\epsilon\bar{\nu}\omega\lambda$  ‘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 ΑΝ 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 μή in μή-οὔ ‘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 ἄρτοι ‘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 τοῦ ‘in this way’  
 (42) ὅς/PDEM τοῦ τοῦ ‘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/ῥμο is regarded as a preposition. The following examples illustrate these principles.

- |      |                  |   |
|------|------------------|---|
| (50) | ετβε/PREP ογ     | ‘for what? why?’                          |
| (51) | εβολ/ADV ῥῥ/PREP | ‘from, out of’ (lit. ‘out in’)            |
| (52) | εχῥ/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) | μμο/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 σ ‘say it!’  
 (59) ρ̄ π σοφ̄ι/N ῃ π νογτε ‘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:  $\rho\iota\nu\omicron\upsilon\epsilon$  ‘to sin’ is considered as single, compound verb (which can still be analyzed morphologically into two units  $\rho+\iota\nu\omicron\upsilon\epsilon$ , 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’,  $\lambda\alpha\alpha\gamma$  ‘something’ and  $\rho\omicron\iota\nu\epsilon$  ‘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)  $\alpha\rho\iota\tau\mu\epsilon\lambda\lambda$ /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’,  $\mu\alpha\nu\omicron\gamma\epsilon$  ‘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  $\tau$ /PPERS ‘he said’  
 (65)  $\mu\alpha\nu\omicron\gamma$ /VBD  $\varsigma$ /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  $\tau$ /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  $\tau\ \mu\alpha\sigma\omega\ \lambda\alpha\tau\epsilon\tau\eta\gamma\tau\iota$  ‘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.