

# Coptic SCRIPTORIUM – Lemmatization Guidelines

Version 1.0.5 / 2017-09-15

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. †, † 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. PERS,

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.	ΑΝΟΚ	ΑΝΟΚ, ΑΝΓ, †, 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.	ΝΤΟϞ	ΝΤΟϞ, Ϟ
3rd sg. fem.	ΝΤΟϚ	ΝΤΟϚ, Ϛ
1st pl.	ΑΝΟΝ	ΑΝΟΝ, ΑΝ, Ν, ΤΝ, CN
2nd pl.	ΝΤΩΤΝ	ΝΤΩΤΝ, ΝΤΕΤΝ, ΤΝ, ΤΗΥΤΝ
3rd pl.	ΝΤΟΟΥ	ΝΤΟΟΥ, Υ, ΟΥ, ΟΕ, ΟΥ

The pronoun lemmas alone therefore primarily give access to search by person (1<sup>st</sup>, 2<sup>nd</sup> ...); 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. ΑΝΓ not ΑΝΟΚ) use the form search norm="ΑΝΓ".

#### *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 πεϞ is πεϞ, the lemma of πα is πα etc., but the lemma of τεϚ and νεϚ is also νεϚ. 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 πεϞ, τεϞ, νοϞ ...), all third person plural possessives (lemma="πεϞ", finds πεϞ, τεϞ and νεϞ) and all third person plural possessives of feminine objects (norm="τεϞ"), 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 -ωϚ, which are lemmatized as such, e.g. ζωλωϚ has the lemma ζωλωϚ.

#### **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 ρβηγϞ, and similarly, possessed forms like τοστ(Ϟ) are lemmatized to the absolute form, i.e. τωρε. 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.  $\zeta\eta\lambda\alpha\varsigma$  and  $\zeta\eta\epsilon-$  ‘(one’s) will’ are lemmatized as  $\zeta\eta\epsilon$ . If only a presuffixal form exists, it is taken as the lemma as well, e.g.  $\eta\lambda\iota\alpha\tau\varsigma$  ‘blessed is...’ has the lemma  $\eta\lambda\iota\alpha\tau$ .

Nouns that have related masculine and feminine forms are considered separate lemmas. For instance, the noun  $\omega\eta\rho\epsilon$  ‘son’ is its own lemma, and the separate noun  $\omega\epsilon\epsilon\rho\epsilon$  ‘daughter’ also has a separate lemma (which is  $\omega\epsilon\epsilon\rho\epsilon$ ). Similarly, Greek words in  $-\omicron\varsigma$  are considered separate from related words in  $-\omicron\eta$ , e.g.  $\pi\omicron\eta\eta\rho\omicron\varsigma$  ‘wicked person’ is its own lemma, and so is the separate  $\pi\omicron\eta\eta\rho\omicron\eta$  ‘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.  $\sigma\omicron\tau\tau\epsilon$  and  $\sigma\epsilon\tau\tau-$  are lemmatized as  $\sigma\omicron\tau\tau$  ‘choose’. The same applies to stative and imperative forms, which are lemmatized to the dictionary entry, e.g.  $\kappa\eta\tau^\dagger$  has  $\kappa\omicron\tau$  as a lemma and  $\lambda\rho\iota$  has  $\epsilon\iota\rho\epsilon$ . Likewise for prenominal forms,  $\sigma\epsilon\tau\tau$  and  $\rho$  are lemmatized as  $\sigma\omicron\tau\tau$  and  $\epsilon\iota\rho\epsilon$ . Compound imperatives receive compound infinitive forms, i.e. for  $\lambda\rho\iota\theta\upsilon\varsigma\iota\alpha\zeta\epsilon$  ‘sacrifice!’, the lemma is  $\rho\theta\upsilon\varsigma\iota\alpha\zeta\epsilon$ .

Note that auxiliaries are not lemmatized to their etymological verbs, i.e. the lemma of the past tense  $\lambda-$  is not  $\epsilon\iota\rho\epsilon$  but  $\lambda$ . Additionally, the negative imperative marker  $\mu\eta\rho$  is lemmatized as  $\mu\eta\rho$  as well, as it is considered to be a form of negation independent from the verb  $\epsilon\iota\rho\epsilon$ . However, the negative imperative of  $\epsilon\iota\rho\epsilon$  itself,  $\mu\eta\omega\rho$  *IS* lemmatized as  $\epsilon\iota\rho\epsilon$  (since it is a morphological imperative of  $\epsilon\iota\rho\epsilon$  itself, and functions as part of its paradigm with the sense ‘to do’).

For fused verb-object forms like  $\eta\tau$  ‘bring me’, see Portmanteau Tags.

### Prepositions

Prepositions are lemmatized to their standard form **before noun phrases**. Therefore the lemma of  $\epsilon-$  and  $\epsilon\rho\omicron-$  is  $\epsilon$ . For preposition forms containing a second person singular feminine pronoun (realized as zero), e.g.  $\epsilon\chi\omega$  ‘on you (fem.)’,  $\eta\sigma\omega$  ‘behind you (fem.)’ etc. see Portmanteau Tags.

### Existential and possessive predicates

The existential predicates are lemmatized as  $\omicron\upsilon\eta$  ‘there is’ and  $\mu\eta$  ‘there isn’t’ (again note that lemmatization does not contain supralinear strokes). Like auxiliaries, the related

possessive predicates are lemmatized using their form before the third person masculine singular: ΟΥΝΤΑ and ΜΗΝΤΑ.

### **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 -ε: in normalized orthography, this is generally present before a nominal subject. The lemma of ΜΑΡΕ- and ΜΑΡ- (jussive) is ΜΑΡΕ, and the lemma of ΩΔΑΝΤ- and ΩΔΑΝΤΕ is ΩΔΑΝΤΕ.

However, the lemmas of auxiliaries that sometimes contain an intermediate pronoun do not contain that pronoun when they occur uninterrupted: the lemmas of ΕΡΩΔΑΝ (conditional) and ΕΡΕ (optative) remain ΕΡΩΔΑΝ and ΕΡΕ. 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 Ν, ΔΝ and ΤΜ are their own lemmas (the form Μ before a labial is also lemmatized as Ν). The negative imperative marker ΜΠΡ is lemmatized as itself (ΜΠΡ), and NOT as ΕΙΡΕ.

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

### **Converters**

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

CCIRC/CFOC:	ΕΡΕ
CREL:	ΕΤΕΡΕ
CPRET:	ΝΕΡΕ

For second person singular feminine ΕΡ/ΕΡΕ (lemma="ΕΡΕ\_ΝΤΟ") see Portmanteau Tags.

### **Inflected modifiers**

Modifiers of the type ΖΩΩ-, ΜΜΙΝΜΜΟ-, ΜΑΥΑΔ-, ΤΗΡ- are lemmatized to their form before the **third person masculine singular** pronoun ς. Thus ΜΜΙΝΜΜΟ- and ΜΜΙΝΜΜΩ- are lemmatized as ΜΜΙΝΜΜΟ. The portmanteau form ΜΜΙΝΜΜΟ (yourself, fem. sg.) is lemmatized ΜΜΙΝΜΜΟ\_ΝΤΟ (see Portmanteau Tags).

## Numerals

Feminine and masculine numerals take the masculine form as the lemma in order to facilitate searches based on the quantity itself (specific searches for either gender can be done using the literal form). For example the lemma of  $\epsilon\alpha\omega\upsilon\epsilon$  ‘seven (fem.)’ is  $\epsilon\alpha\omega\upsilon$  ‘seven (masc.)’. Note that compound numbers receive a complex lemma, therefore the lemma of  $\chi\omicron\upsilon\gamma\omega\tau\epsilon\alpha\omega\upsilon\epsilon$  ‘twenty-seven (fem.)’ is  $\chi\omicron\upsilon\gamma\omega\tau\epsilon\alpha\omega\upsilon$  ‘twenty-seven (masc.)’ (the individual parts can still be annotated at the morph level, which is not lemmatized).

## Portmanteau Tags

Some fused items receive a so-called portmanteau tag representing two categories at once. For example, the form  $\epsilon\varphi\omega\lambda\eta$  is considered to contain a conditional auxiliary and a subject pronoun: `pos="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  $\epsilon\iota\omega\lambda\eta$ ,  $\epsilon\sigma\omega\lambda\eta$  and  $\epsilon\varphi\omega\lambda\eta$  are  $\epsilon\varphi\omega\lambda\eta\_a\lambda\omicron\kappa$ ,  $\epsilon\varphi\omega\lambda\eta\_n\tau\omicron\varsigma$  and  $\epsilon\varphi\omega\lambda\eta\_n\tau\omicron\upsilon$  respectively. The lemma of  $\epsilon\varphi\omega\lambda\eta$  remains  $\epsilon\varphi\omega\lambda\eta$  (`pos="ACOND"`), unless it contains a second person feminine singular subject, in which case the lemma is  $\epsilon\varphi\omega\lambda\eta\_n\tau\omicron$  according to the rule above.

For the past tense second person singular feminine form  $\lambda\upsilon$  the lemma is similarly  $\lambda\_n\tau\omicron$  (`pos="APST_PPERS"`). The form  $\mu\mu\iota\eta\mu\mu\iota$  (yourself, fem. sg.) is identical to the base of other personal forms, but is lemmatized  $\mu\mu\iota\eta\mu\mu\iota\_n\tau\omicron$  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  $\epsilon\chi\eta\_n\tau\omicron$ ,  $\eta\sigma\omega$  ‘behind you (fem.)’ has  $\eta\sigma\alpha\_n\tau\omicron$  etc.

For circumstantial or focalizing converter + second person feminine singular, the lemma  $\epsilon\varphi\epsilon\_n\tau\omicron$  is used (and similarly preterit  $\eta\epsilon\varphi\epsilon\_n\tau\omicron$  and relative  $\epsilon\tau\epsilon\varphi\epsilon\_n\tau\omicron$ ).

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:  $\epsilon\iota\eta\epsilon\_a\lambda\omicron\kappa$ .

## Confusing cases

### Nouns with variant spellings

Lemmatization is the highest level of lexical abstraction, and as such should unify variants, even if the underlying forms are alternative ‘canonical’ spellings (e.g. listed in Crum’s dictionary as subentries). Spelling variants that are canonical should be retained on the norm level in order to make them findable as well – this is distinguished from the

orig level, which is the least abstract, and may additionally contain non-canonical spelling variation, as well as diacritics.

Examples:

<b>orig form</b>	<b>norm</b>	<b>lemma</b>	<b>(in bound group)</b>	<b>translation</b>
οογ	οογ	ζοογ	ῆποογ	today
ῆμν	μν	μν	ῆμνρωμε	there isn't

In other words, because the form ποογ is a standard spelling for ηζοογ in ‘today’, we do not normalize with a hori, allowing users to find this variant in norm while abstracting away from potential diacritics (including damage to letters, etc.), which would be present in orig. Users wishing to find all cases of the word for ‘day’ can still do so by searching for the lemma ζοογ.