

Coptic Scriptorium – Entity Annotation Guidelines

Amir Zeldes

Georgetown University

Version: 1.0.0_2020.05.20

1 Introduction

1.1 Preamble

Entity annotation concerns the annotation of **referring expressions** in a text, i.e. spans of text that refer to things in the world, and their classification into **entity types**. The purpose of entity annotation in Coptic Scriptorium is to facilitate searches which include specific entity types (e.g. finding a certain epithet using linguistic annotations, such as ⲟⲩⲁⲁⲃ ‘holy’, but only when applied to a PERSON), to inventorize entities (find all cases of e.g. places mentioned in the Apophthegmata Patrum), and to function as a gateway for entity linking, enabling searches for specific persons (“John the Baptist”), regardless of the exact expression used to mention them. The latter task of entity linking is left outside of the scope of the current guidelines.

Entity annotation can be applied to three types of referring expressions:

- Named entities, which are headed by a proper noun (e.g. “Apa Papnoute”)
- Non-named entities, headed by a common noun (e.g. “the angel”)
- Pronouns – these are currently not annotated by our schema (e.g. “she” is a person)

1.2 Referring expressions

Almost all nouns and proper nouns correspond to referring expressions, with the exception of non-referring nouns, such as:

- ⲁⲗⲉ ⲣⲁⲧ.. - “stand, set foot” - does not actually refer to the foot of a person

- ζῆν οὐ μὲν - “truly” - does not actually introduce a referenceable ‘truth’

One test for referentiality is whether a pronominal or nominal subsequent mention is possible/plausible. For example, the following sounds odd:

- ?? ἀφάρξατο τὴν ἀπὸ τοῦ ποδὸς ... “he stood on foot, and this foot...”

For more examples, also see the section “Non-referring expressions” below.

2 Entity Types

We distinguish 10 entity types:

- ABSTRACT - intangible entities not covered by other classes (incl. ideas, emotions)
- ANIMAL – dog, fish, ...
- EVENT – an occurrence, e.g. “the death of the king”, “the arrival of a monk”
- OBJECT – concrete inanimate object not belonging to other categories
- ORGANIZATION – organized body of people, e.g. ἑκκλησία, τὸ στρατεῖα
- PERSON – references to humans, loose groups of humans (ὄχλος ‘crowd’), deities
- PLACE – towns, countries, but also ad-hoc places (behind the house, outside)
- SUBSTANCE – mass noun indicating a material, e.g. sand, water, wine
- TIME – date terms, durations like ‘year’, ‘day’, terms like ‘moment’

3 Markable selection guidelines

3.1 Appositions

Repeated mentions of the same entity in apposition are considered a single span, and do not contain more mentions of the same entity:

- [Ἰωάννης ὁ βαπτιστής] “[John the Baptist]”
- [ὁ πρὸς Ζήνων] “[King Zeno]”
- [πᾶν τὸ ἐν ὁυδὶ κατὰ σμὸς τοῦ ἀγίου κυροῦ πᾶν ἅπαντα ...] “[Our Holy One in every way, Apa Cyrus, who has ...]”

Although outwardly similar, appositions must be distinguished from dislocations, in which a pronominal subject or object is repeated separately. For personal pronouns, the pronoun is simply left out of the nominal span:

- [περ ειωτ] ϣ ναγ ερο ου - “[his father], he sees them”
- ϣ ναγ ερο ου νσι [περ ειωτ] - “he sees them, that is [his father]”

If the pronoun is a substitutive demonstrative (παι, ται, ναι), then two spans are annotated:

- [περ ειωτ] [παι] ναγ ερο ου - “[his father], [this one] sees them”
- [παι] ναγ ερο ου νσι [περ ειωτ] - “[this one] sees them, [that is his father]”

But note that it is also possible for a substitutive demonstrative to stand in true apposition to a noun without dislocation, in which case a single span is annotated as for any apposition:

- α ι ναγ ε [περ ειωτ , παι ετ μεριτ ς] - “I saw [their father, the one who loves her]”

See the UD Coptic guidelines for more information on identifying dislocation vs. apposition.

3.2 Expanded Relative Constructions

The relative construction expanding an article is annotated as an entity:

- [π ετ ου σωτμ ερο ϣ] “the one they listened to” (person)

However, if the π is tagged as a copula, that part of the construction is not part of the entity span, since it is part of a predication. In these instances, we view the predicate noun phrase as an entity, and the relative clause as a subject clause (compare the Universal Dependency annotation guidelines):

- [π νογτε] π/COP εντ α ϣ αγζανε “It is God who made them grow”

In this example, “God” receives a span, but “who made them grow” is considered a subject clause (i.e. ‘who made them grow is God’), which is not nominal and hence not annotated. Note that according to the tagging guidelines, the second π should be tagged as COP and lemmatized πε in this sentence.

3.3 Possessive Constructions

The possessive article construction, e.g. *τα παγλολο*, forms two spans, as follows, with the entity type being decided based on meaning:

- [*τα παγλολο*]_{PERSON}_{PERSON} - the ones (=people) belonging to Paul

But note that regular possessive articles are not annotated with spans, just as other pronouns are not annotated:

- [*περ η*]_{PLACE} - his house

3.4 Interrupted spans

Entity expressions interrupted e.g. by a copula or particle are spanned to **contain** the copula or particle. For example, the following span includes the intervening copula:

- [*περ αποστολο* ***νε*** *ετογαδ*]_{PERSON} ‘it is his holy apostles’ (literally [*his apostles are which holy*], with intruding ‘are’)

Similarly:

- [*πτοογ δε η ροογ*]_{TIME} - but four days (lit. ‘[four but days]’)
- [*ογ ωβηρ · ρωω κ ον νε π νογτε*]_{PERSON} - but also for your part a friend of God

Non-adjacent relative clauses are included, **unless the interruption contains the verb controlling the head noun** (this prevents some possibly very long ‘hermeneutical’ relatives inside mentions):

- [*ρωμε νη ον ετ ωτμ*] - and also [any man who hears] (note the interruption ‘ον’, and inclusion of the relative clause)

But do not include a clause past the verb controlling the head of the span:

- *ερωαν* [*τ βαωορ*]_{ANIMAL} *αωκακ εβολ αν* ***ετε ντοκ πε*** ... - it is not when [the fox] barks, which is you, ... (postponed hermeneutic relative clause in bold is not included, because it appears after the verb *αωκακ*, which controls ‘fox’ as a subject)

In this case the interruption by the verb *αωκακ* ‘bark’ which is the predicate of ‘fox’ triggers the guideline to omit the relative clause. Otherwise, the mention could potentially cover the entire

clause, in this case: τ βαϥορ αϥκακ εβολ αν ετε ντοκ πε π ρμζαλ η π μαμμωνας ρν ρεν ρροϥ
ε γ οϥ

3.5 Groups and other quantity constructions

Semantically ‘empty’ heads such as quantity nouns (compare English ‘a number of people’, which is not both ‘a number’ and ‘people’; similarly ‘a lot of’, ‘the majority of’ etc.) are only given one span, for example:

- [ραζ η κοτ]_{TIME} – a lot of times

Groups of entities are generally interpreted as the entity type of their constituents, for example, a herd of animals is of the type animal:

- [οϥ αρελη η ϥοϥ]_{ANIMAL} - a herd of buffaloes

Note that there is no nested entity for 'buffalo' in this case, since there is no distinct entity being mentioned (the herd consists exactly of all buffaloes being discussed). This is different in cases where the nested entity is not identical in reference, e.g. ‘[the houses of [the city]]’, where ‘the city’ can be said to contain more than just the houses.

An exception to the guideline that groups are classified as their constituent type is cases of people who form an organization, e.g. ϥυναγωγη, στρατευμα etc. are 'organization', not 'person'.

3.6 No reference inside compounds

In morphologically complex items containing a verb inside a larger token, that noun cannot be annotated:

- α ϥ ριβαπτισμα - he received-baptism

In this case baptism cannot be annotated as an entity, since it's part of an incorporated verb ‘to baptize’, and receives the part of speech V in Coptic Scriptorium guidelines.

3.7 Coordination

Do not mark coordinate entities in addition to their constituents:

- [ιωζαννης] ην [αντωνιος]

In this case we do not also annotate [ΙΩΡΑΝΝΗΣ ΜΗ ΑΝΤΩΝΙΟΣ] as a third mentioned entity.

3.8 Container and substance

Container and substance form two entities, for example:

- [ΟΥ ΠΥΓΗ Μ [ΜΟΟΥ]_{SUBSTANCE}]_{PLACE} - a fountain of water

The fountain can be a PLACE or OBJECT in context, but the water is SUBSTANCE, and both can be referred to separately later on.

3.9 Distributive entities

Repeated distributive noun constructions are interpreted as single entity mentions:

- [Π ΟΥΑ Π ΟΥΑ]_{PERSON} - ‘one by one’, ‘each man’

The rationale is that these are like a plural reference, rather than two mentions of individuals (in this case there can be more than two people, and they do not map neatly onto the two numerals).

4 Entity classification guidelines

4.1 Body Parts

Most body parts are marked as objects, since they are tangible:

- [ΟΥ ΓΙΧ]_{OBJECT} – “a hand”
- [ΠΕΦ ΒΑΛ]_{OBJECT} – “his eye”

However some referential body parts are usually considered abstract, notably ΖΗΤ ‘heart’, which usually refers to one’s spirit, emotions, etc., and not the physical organ:

- † ΝΑ ΤΡΕ [ΠΟΥ ΖΗΤ]_{ABSTRACT} ΜΚΑΞ – “I will make your [heart] suffer”

Other uses of body parts may be totally figurative or idiomatic (i.e. not referring to anything), in which case they are not annotated – see ‘Non-Referring Expressions’ above.

4.2 Peoples and demonyms

Pluralized demonyms indicating members of a people are labeled person:

- [ἡ ῥελλην]_{PERSON}

However peoples mentioned as a people (not as a group of individuals) are labeled organization:

- [ἡ ἐκ λαοῦ ἰσραὴλ]_{ORGANIZATION}

These cases are usually singular and involve a named people. This guideline does not apply to ad-hoc groups of people who do not form an organized entity, e.g. ὄχλος ‘crowd’ is still annotated as PERSON.

4.3 Multiple entity types in one chain

It is possible for an entity to be referred to in multiple ways, underscoring different aspects of the entity. In such cases, where it seems clear that the entity types are distinct (e.g. due to metonymy, metaphorical extensions, etc.), it is possible to have different entity types

- [ἡ χριστός]_{PERSON} ἡ [τῆς ἐκτῆς] _{ABSTRACT} - ‘this foundation is Christ’ (where ‘foundation’ refers to ‘Christ’ metaphorically)

Mark each entity with its own type, i.e., Christ as PERSON and Foundation as ABSTRACT.

5 Non-referring expressions

5.1 Interrogatives

No annotations are needed for plain interrogatives (τίς, τίς), but complex interrogatives including a noun are annotated:

- [τίς γάρ ὁ ἄνθρωπος] _{PERSON} οὗτος οὖν ὁ [ὁ ἀνθρώπος] _{PERSON} “[what human] is he who knows [those things which are of humans]]?”
- [τίς ἡ πειρασμός] _{EVENT} “[what kind of trial]?”

5.2 Common figurative and other fixed expressions to ignore

The following are considered idiomatic or functional expressions, in which the constituent nouns are not construed as referential and no annotation is needed:

- ἵστημι ποῦ - ‘stand, set foot’ - “foot” is not an entity mention
- βοηθεῖν ποῦ - ‘help, give a hand’
- ἐπὶ τὴν γῆν - ‘down’, lit. ‘to the ground’

- Ν ΟΥ ΖΟΥΟ - ‘more so’
- (ΖΩΒ) Ν ΓΙΧ – ‘handywork’ - the whole phrase (handywork) is ABSTRACT or OBJECT in context, but ‘hand’ is not a referent
- Ρ ΖΝΑ Υ – ‘want, do one’s will’ - the word ΖΝΑ / ΖΝΕ ‘will’ is figurative, as this is a fixed expression for ‘desire’
- Ε Π ΤΗΡΥ meaning ‘at all’ is not referential
- ΖΑ ΕΟΥ - ‘glorious’ - the ΕΟΥ is not referential, as the expression is only used adjectivally
- Ν ΟΥ ΚΟΥΙ - ‘a little’ (manner adverbial; note ΚΟΥΙ *can* be referential if referring to a person or thing, e.g. ΝΤΚ [ΟΥ ΚΟΥΙ]_{PERSON} ‘you are a little one’)
- Ν Τ ΖΕ - meaning ‘like’
- Ν ΦΟΡΠ - ‘first(ly)’
- ΒΟΜ - meaning ‘capable’ in constructions like ΝΤΚ ΒΟΜ ΔΝ ‘you are not capable’, but ΒΟΜ can be referential in, e.g. [Τ ΒΟΜ Μ [Π ΝΟΥΤΕ]]_{ABSTRACT}
- Ν ΟΥΩΤ - ‘together’
- ΖΙ ΟΥ ΣΟΠ - ‘at once’
- Μ ΜΗΝΕ - ‘daily’