

Principles of Syriac Lemmatisation

Summary version

Nicolas Atas (Freie Universität Berlin) – August 2021

Lemmatisation is the operation identifying the lemma (headword) of each word-form attested in a given text. Each lemma is also accompanied by its part-of-speech tag (POS-tag) indicating its morphosyntactic category (noun, verb, adjective, pronoun, etc.), and, for Semitic languages as Arabic and Syriac, its root, and its origin. This paper describes these elements in a summarizing way¹. For more information, see the website of the GREgORI Project (including complete bibliography)².

1. Word-forms

Syriac words are divided into two categories, monolexical word-forms and polylexical word-forms. Our system of analysis takes into account each of the constitutive lexical elements of these forms, whether they are monolexical or polylexical.

1. Monolexical word-forms which are constituted by a single word and can be directly related to a single lemma:

- a. ܐܘܪܝܬܐ (<ܐܘܪܝܬܐܐܘܪܝܬܐ “letter” Sokoloff, p. 9)
- b. ܡܢܩܕܡܐ (<ܡܢܩܕܡܐ “before” Sokoloff, p. 1318)
- c. ܐܠܘܗܐ (<ܐܠܘܗܐܐܠܘܗܐ “God” Sokoloff, p. 47)
- d. ܐܢܝ (<ܐܢܝܐܢܝ “I” Sokoloff, p. 58)
- e. ܡܢܐ (<ܡܢܐܡܢܐ “what” Sokoloff, p. 778)
- f. ܡܘܨܥܐ (<ܡܘܨܥܐܡܘܨܥܐ “to find” Sokoloff, p. 1556)
- g. ܕܘܪܝܬܐ (<ܕܘܪܝܬܐܕܘܪܝܬܐ “to desire” Sokoloff, p. 1435)
- h. ܡܘܨܦܐ (< 2 ܡܘܨܦܐܡܘܨܦܐ “to understand” Sokoloff, p. 1010)
- i. ܝܫܘܥܐ (<ܝܫܘܥܐܝܫܘܥܐ “Jesus” Payne Smith, col. 1638)

2. Polylexical word-forms which first need to be divided into multiple *word-forms* before they can be related to a lemma³:

- a. ܡܘܨܥܐ: ܡܘܨܥܐܡܘܨܥܐ (<ܡܘܨܥܐ “soul” Sokoloff, p. 938 and ܡܘܨܥܐ “he” Sokoloff, p. 333)
- b. ܡܘܨܥܐ: ܡܘܨܥܐܡܘܨܥܐ (<ܡܘܨܥܐ “and” Sokoloff, p. 357 and ܡܘܨܥܐ “henceforth” Sokoloff, p. 758)
- c. ܡܘܨܥܐ: ܡܘܨܥܐܡܘܨܥܐ (<ܡܘܨܥܐ “and” Sokoloff, p. 357 and ܡܘܨܥܐ “to pray” Sokoloff, p. 1288)
- d. ܡܘܨܥܐ: ܡܘܨܥܐܡܘܨܥܐ (<ܡܘܨܥܐ “in order that” Sokoloff, p. 268 and ܡܘܨܥܐ “to find” Sokoloff, p. 1556)
- e. ܡܘܨܥܐ: ܡܘܨܥܐܡܘܨܥܐ (<ܡܘܨܥܐ “to, for” Sokoloff, p. 665 and ܡܘܨܥܐ “labor” Sokoloff, p. 1112)
- f. ܡܘܨܥܐ: ܡܘܨܥܐܡܘܨܥܐ (<ܡܘܨܥܐ “and” Sokoloff, p. 357; ܡܘܨܥܐ “love” Sokoloff, p. 419 and ܡܘܨܥܐ “he” Sokoloff, p. 333)
- g. ܡܘܨܥܐ: ܡܘܨܥܐܡܘܨܥܐ (<ܡܘܨܥܐ “and” Sokoloff, p. 357; ܡܘܨܥܐ “in, through” Sokoloff, p. 114; ܡܘܨܥܐ “suffering” Sokoloff, p. 497 and ܡܘܨܥܐ “we” Sokoloff, p. 472)
- h. ܡܘܨܥܐ: ܡܘܨܥܐܡܘܨܥܐ (<ܡܘܨܥܐ “and” Sokoloff, p. 357; ܡܘܨܥܐ “to, for” Sokoloff, p. 665; ܡܘܨܥܐ “glory” Sokoloff, p. 1518 and ܡܘܨܥܐ “he” Sokoloff, p. 333)
- i. ܡܘܨܥܐ: ܡܘܨܥܐܡܘܨܥܐ (<ܡܘܨܥܐ “and” Sokoloff, p. 357; ܡܘܨܥܐ “of” Sokoloff, p. 268; ܡܘܨܥܐ “negligence” Sokoloff, p. 720 and ܡܘܨܥܐ “we” Sokoloff, p. 472)

¹ Thanks are due to all our collaborators in the field of Syriac lemmatisation for the GREgORI Project (alphabetical order): Naima Afif, Yury Arzhanov, Philip Michael Forness, Jean-Claude Haelewyck, Bastien Kindt, Manhal Makhoul, David Phillips, Marcel Pirard, Andrea Barbara Schmidt, Antonio Stefano Sembianti and Guido Venturini. A special gratitude is also due to Andrea Barbara Schmidt and Philip Michael Forness, who reviewed these lines and their English translation.

² See <https://uclouvain.be/fr/instituts-recherche/incal/ciol/gregori-project.html>.

³ Polylexical word-forms are indicated in the corresponding column by the letter K (from the Greek “crasis”, contracted forms, appellation also used for prefixed and suffixed forms) and lexical elements included in polylexical word-forms are divided by means of the sign “@”.

2. Lemmas

Lemmas represent words as they appear in dictionaries. Our system of Syriac lemmatisation is based on Sokoloff's *Syriac lexicon*⁴ used as “reference dictionary”:

1. Lemmas correspond to the headwords found in Sokoloff's dictionary with eastern vocalization, all diacritical marks (including spirantization and linea occultans), in the emphatic state⁵:
 - a. Form ܐܘܠܡܐ: lemma ܐܘܠܡܐ (“love” Sokoloff, p. 419)
 - b. Form ܐܘܠܡܐܐ: lemma ܐܘܠܡܐܐ (“new” Sokoloff, p. 418)
2. Lemmas of proper names (anthroponyms and toponyms) follow the spelling found in Payne Smith's *Thesaurus*⁶, without vocalization:
 - a. Form ܕܘܥܘܫܐ: lemma ܕܘܥܘܫܐ (“Jesus” Payne Smith, col. 1638)
 - b. Form ܕܘܥܘܫܐܐ: lemma ܕܘܥܘܫܐܐ (“Dadišo” Payne Smith, col. 824)
3. Homographic lemmas are distinguished by a number conforming to the usage of Sokoloff's dictionary:
 - a. Form ܡܘܨܦܐ: lemma “1” ܡܘܨܦܐ (“to hope” Sokoloff, p. 964)
 - b. Form ܡܘܨܦܐܐ: lemma “2” ܡܘܨܦܐܐ (“to announce” Sokoloff, p. 965)
4. Word unrecorded neither in the reference dictionary nor in the database of the GREgORI project is recorded with lemma as intuitively expected by linguists:
 - a. Form ܕܘܥܘܫܐܐܘܠܡܐܐ: lemma ܕܘܥܘܫܐܐܘܠܡܐܐ “being recompensed”
 - b. Form ܕܘܥܘܫܐܐܘܠܡܐܐܐ: lemma ܕܘܥܘܫܐܐܘܠܡܐܐܐ “prone to disease”

3. POS-tags

Each lemma receives a tag that indicates the part of speech, the morphosyntactic category, to which it belongs. This tag characterizes the lemma at the lexical level in general. It does not describe in any way the particular use of the lemma in a specific context. POS-tags used in our system of analysis are as follows:

- ADJ**: adjective (ܡܘܨܦܐ “lacking” Sokoloff, p. 1024)
- ADV**: adverb (ܐܘܕܘܐ “again” Sokoloff, p. 1626)
- CARD**: cardinal number (ܒܝܕ “one” Sokoloff, p. 413)
- INTJ**: interjection (ܘܘܐ “woe!” Sokoloff, p. 357)
- NAME**: proper noun
- NAME_ant**: anthroponym ܕܘܥܘܫܐ “Adam” (Payne Smith. col. 38)
- NAME_top**: toponym ܡܘܨܦܐ “Egypt” (Payne Smith. col. 2196)
- NOUN**: noun (ܐܘܠܡܐ “brother” Sokoloff, p. 25)
- ORD**: ordinal number (ܦܘܪܝܫܐ “first” Sokoloff, p. 1319)
- PART**: particle⁷ (ܘܘܐ “and” Sokoloff, p. 357)
- PRO_dem**: demonstrative pronoun (ܐܘܠܡܐ “this one” Sokoloff, p. 346)
- PRO_ind**: indefinite pronoun (ܡܘܨܦܐ “something” Sokoloff, p. 715)
- PRO_int**: interrogative pronoun (ܐܘܠܡܐ “how” Sokoloff, p. 34)
- PRO_pers**: personal pronoun (ܐܘܠܡܐ “I” Sokoloff, p. 58)
- V1–V33**: verbs⁸
- V1**: ܡܘܨܦܐ “to confine” (Sokoloff, p. 411)
- V3**: ܡܘܨܦܐܐ “to announce” (Sokoloff, p. 965)
- V5**: ܐܘܠܡܐܐ “to terrify” (Sokoloff, p. 374)
- ...

⁴ M. SOKOLOFF, *A Syriac Lexicon: A Translation from the Latin, Correction, Expansion, and Update of C. Brockelmann's Lexicon Syriacum*, Piscataway, 2009 (2d corrected print, 2012).

⁵ For the sake of readability, we will only present from this point forward monolexical occurrences. Passive participles in the absolute state in Sokoloff will still receive a lemma in the emphatic state.

⁶ R. PAYNE SMITH, *Thesaurus syriacus*, Oxford, 1871-1901 (reprint Hildesheim, 2006).

⁷ Particles are tagged as conjunction and preposition in Sokoloff's *Syriac Lexicon*.

⁸ Complete table of the verbal POS-tags can be found at the following address: https://www.gregoriproject.com/pdf/POS_SYC.pdf.

4. Roots

At the level of the root, two categories are distinguished.

1. **Base Nouns** (primary nouns)⁹ (vocalized, if the vocalization is attested):
 - a. Lemmas: ܠܫܢܐ, ܠܫܢܐܘܬܐ and ܠܫܢܐܘܬܐܘܬܐ: root ܠܫܢܐ (Sokoloff, p. 1)
 - b. Lemmas: ܠܫܢܐܘܬܐ, ܠܫܢܐ: root ܠܫܢܐܘܬܐ (Sokoloff, p. 965)
2. **Consonantal Roots** (not vocalized), linked to a base verb from which nouns and adjectives are derived:
 - a. Lemmas: ܠܫܢܐܘܬܐ, ܠܫܢܐܘܬܐ, ܠܫܢܐܘܬܐܘܬܐܘܬܐ, ܠܫܢܐܘܬܐܘܬܐ, ܠܫܢܐܘܬܐܘܬܐ: root ܠܫܢܐ (Sokoloff, p. 374)
 - b. Lemmas: ܠܫܢܐܘܬܐ, ܠܫܢܐܘܬܐ, ܠܫܢܐܘܬܐܘܬܐ: root ܠܫܢܐܘܬܐ (Sokoloff, p. 1352)

Particular cases:

1. Homographic roots are distinguished by a number conforming to the usage of Sokoloff’s dictionary:
 - a. Lemma 1 ܠܫܢܐ: root 1 ܠܫܢܐ
 lemma 2 ܠܫܢܐ: root 2 ܠܫܢܐ
 lemma 3 ܠܫܢܐ: root 3 ܠܫܢܐ
 lemma 4 ܠܫܢܐ: root 4 ܠܫܢܐ
 (as Sokoloff, pp. 1133-1134)
 - b. Lemmas 1 ܠܫܢܐܘܬܐ, 1 ܠܫܢܐܘܬܐܘܬܐ, 1 ܠܫܢܐܘܬܐܘܬܐܘܬܐ, 1 ܠܫܢܐܘܬܐܘܬܐܘܬܐ and 1 ܠܫܢܐܘܬܐܘܬܐܘܬܐ: root 1 ܠܫܢܐ;
 lemmas 2 ܠܫܢܐܘܬܐܘܬܐ and 2 ܠܫܢܐܘܬܐܘܬܐܘܬܐܘܬܐ: root 2 ܠܫܢܐ
 (as Sokoloff, pp. 1044-1045)
2. If the word is not recorded in reference dictionaries:
 - a. If the root (consonantal or base noun) is known, this root is used
 - i. Lemma ܠܫܢܐܘܬܐܘܬܐܘܬܐ “being recompensed”: root 3 ܠܫܢܐ
 - ii. Lemma ܠܫܢܐܘܬܐܘܬܐ “prone to sickness”: root ܠܫܢܐ
 - b. If the root remains unknown, lemma is used as root
 - i. Lemma ܠܫܢܐܘܬܐܘܬܐܘܬܐ: root (base noun) ܠܫܢܐܘܬܐܘܬܐܘܬܐܘܬܐ
 - ii. Lemma ܠܫܢܐܘܬܐܘܬܐܘܬܐ: root (base noun) ܠܫܢܐܘܬܐܘܬܐܘܬܐܘܬܐ
 - iii. Lemma ܠܫܢܐܘܬܐܘܬܐܘܬܐܘܬܐ: root (base noun) ܠܫܢܐܘܬܐܘܬܐܘܬܐܘܬܐܘܬܐ

5. Origin

Each root receives an indication specifying if the word is borrowed or not from another language than Syriac. In our database, Syriac words are characterized by the sign “0” and non-Syriac word by the sign “-1” (e.g. ܠܫܢܐ, from Greek ܠܫܢܐ). In our concordances in PDF format, non-Syriac roots are therefore automatically characterized by the sign “L” (for “loanword”). Proper names are not taken into account and are conventionally characterized by the “0” sign.

⁹ In this case, it is the noun or the adjective that is original. The majority of base nouns falls into the following categories: family (father, brother, etc.), body parts (heart, foot, etc.), animals (bull, rabbit, etc.), plants (herbs, garlic, etc.), stones and metals (clay, gold, etc.), temporal concepts (year, eternity, etc.), geographical concepts (earth, rivers, etc.), architectural concepts (dike, roof, etc.), tools (basket, vase, etc.) and foreign words (palace, paradise, etc.). They can themselves be at the origin of verbs, which are qualified as *denominative verbs*, generally used in a derived *form*. Taking into account the base nouns obviates the need to create artificial verbal roots that are supposed to account for them (see B. KINDT, J.-Cl. HAELEWYCK, A.B. SCHMIDT, N. ATAS, *La concordance bilingue grecque-syriaque des Discours de Grégoire de Nazianze*, in *BABELAO*, 7 (2018), p. 63-64 ; document freely available [here](#)).

Appendix: Keyboard Shortcuts

	Belgian Keyboard	German Keyboard	English Keyboard
Vowels			
a ܐ (ܐ)	SHIFT + W	SHIFT + Y	SHIFT + Z
ā ܐ̄ (ܐ̄)	SHIFT + X	SHIFT + X	SHIFT + X
e ܐ (ܐ)	SHIFT + C	SHIFT + C	SHIFT + C
ī/ū ܐ̄ (ܐ̄/ܐ̄)	SHIFT + V	SHIFT + V	SHIFT + V
o ܐ (ܐ)	SHIFT + B	SHIFT + B	SHIFT + B
ē ܐ̄ (ܐ̄)	SHIFT + N	SHIFT + N	SHIFT + N
Spirantisation			
<i>Rukkākhā</i> ܐ (ܐ)	SHIFT + ,	SHIFT + M	SHIFT + M
<i>Quššāyā</i> ܐ (ܐ)	SHIFT + U	SHIFT + U	SHIFT + U
Lengthening Line <i>(kashīda)</i> ¹⁰ : - (ܐ̄)			
	SHIFT + J	SHIFT + J	SHIFT + J
Plural Marker <i>(syāmē)</i> : .. (ܐ̄)			
	SHIFT + I	SHIFT + I	SHIFT + I
Linea occultans (<i>mbaṭṭlānā</i>):			
Below - (ܐ̄)	SHIFT + L	SHIFT + L	SHIFT + L
Above - (ܐ̄)	SHIFT + O	SHIFT + O	SHIFT + O
Syriac abbreviation mark: ܐ̄	² (before the word)	^ (before the word)	` (before the word)

¹⁰ In order to distinguish single letter particles ܐ, ܐ̄, ܐ̄, ܐ̄ from numbers, particles are followed by the kashida ܐ̄, ܐ̄, ܐ̄, ܐ̄. In accordance with Sokoloff's dictionary (p. 295), we add also the kashida after the possessive particle -ܐ̄.