

RETRO-DIGITIZATION IN GREEK DIALECTOLOGY AND LEXICOGRAPHY: CHALLENGES OF MORPHO- PHONETIC REPRESENTATION OF THE CAPPADOCIAN DIALECT

Io Manolossou, Athanasios Karasimos & Georgia Katsouda
Academy of Athens

Στόχος της παρούσας μελέτης είναι η παρουσίαση της πρώτης ψηφιακής τυποποιημένης κωδικοποίησης για την σύνταξη ενός ιεραρχικά σχολιασμένου ψηφιακού λεξικού της Καππαδοκικής διαλέκτου. Δεδομένης της ανάγκης για βαθύτερη κατανόηση και αναπαράσταση των διαλεκτικών δεδομένων στα διαφορετικά επίπεδα γλωσσικής ανάλυσης, αλλά και της έλλειψης ενός προτύπου κωδικοποίησης, έχει δημιουργηθεί ένα σχήμα/πρότυπο κωδικοποίησης για πολυεπίπεδη γλωσσική επισημείωση και λεξικογραφική αναπαράσταση. Για τη δημιουργία αυτού του προτύπου αντιμετωπίστηκαν τα ακόλουθα ζητήματα: α.) σύστημα φωνητικής μεταγραφής, τόσο στο IPA όσο και στο ελληνικό αλφάβητο (χρήση ειδικών γραφημάτων, συμμόρφωση στο πρότυπο Unicode, συμβατότητα με προηγούμενες προσεγγίσεις μεταγραφής, φιλικό προς το χρήστη), β.) σύστημα γεωγραφικού προσδιορισμού και γεωθεσίας (τυποποίηση τοπωνυμίων βάσει εξαντλητικής έρευνας των πηγών και αξιολόγηση των παραλλαγών, εναλλακτική υλοποίηση ανάλογα με το μέσο [έντυπο vs ψηφιακό], σύνδεση με ψηφιακό χάρτη και τοπωνύμια από γεωπληροφοριακά συστήματα), και γ.) σύστημα ταξινόμησης και παρουσίασης παραδειγμάτων (βάσει διαφόρων τύπων παραδειγμάτων: αποσπάσματα απομαγνητοφωνημένων προφορικών συνεντεύξεων, παλαιότερες γραπτές πηγές, απλή αφήγηση vs. φράσεις vs. παροιμίες vs. άσματα) και δ.) ένα συνολικό σύστημα παρουσίασης δεδομένων (συμβατότητα με τη μορφή μεγάλων διαλεκτικών λεξικών όπως το Ιστορικό Λεξικό της Νέας Ελληνικής (IANE), το Ιστορικό Λεξικόν των ελληνικών ιδιωμάτων της Κάτω Ιταλίας (IAEIKI), το Ιστορικό Λεξικόν της Ποντιακής διαλέκτου (ΙΠΠΔ) και το Λεξικό Τσακωνικής Διαλέκτου). Το προτεινόμενο πρότυπο συγκρίνεται με αντίστοιχες διεθνείς προσπάθειες στη διαλεκτική ηλεκτρονική λεξικογραφία (π. χ. την ηλεκτρονική έκδοση του Innsbruck EDD Online 3.0), καθώς και με προηγούμενα ελληνικά έργα μικρότερης κλίμακας (π. χ. το Τριδιαλεκτικό Λεξικό).

Keywords: *Cappadocian dialects, dialectology, lexicography, historical dictionary, digital humanities*

1. Introduction

In this paper, we present various aspects of a major on-going digital project funded by the Hellenic Foundation for Research and Innovation (ΕΛΙΔΕΚ). The project, entitled DicaDland (Digitizing the Cappadocian Dialectal Landscape), runs from

2018 until 2021 and aims to produce a full documentation of the Cappadocian dialect, including the varieties of Pharasa and Silli, by exploiting all the available sources (older, written and new, oral) and the latest advances in the domain of digital humanities. It constitutes an intersection of a humanities/social sciences discipline (linguistics) and informatics, as it involves both the examination of how digital tools can be applied to specific domains such as dialectal Lexicography and Modern Greek dialectology, and how these, as domains of application, can influence the development of information technology.

The final products of the project will be (or rather, are, since we are more than half-way-through by now) two state-of-the-art major reference works, namely an interactive electronic dialectal atlas (which will constitute the first such effort in the domain of Greek Linguistics — you'll hear about it later on) and a comprehensive historical dictionary of the Cappadocian dialects (again lacking until now — you'll hear all about it now).

The project¹ is hosted by the Laboratory of Modern Greek Dialects of the University of Patras and co-ordinated by Dimitra Melissaropoulou, Associate Professor at the Aristotle University of Thessaloniki. The second Principal investigator is Athanasios Karasimos, from the Academy of Athens and the Centre for the Greek language, and the research team is complemented by a host of specialists in dialectology, including Io Manolessou and Georgia Katsouda from the Academy of Athens, as well as Angela Ralli, Mark Janse, Petros Karatsareas, Metin Bağrıaçık, Symeon Tsolakidis, Christos Papanagiotou, Georgios Chairetakis, Stavros Bompolas, and Dimitris Papazachariou.

The importance of the Cappadocian dialect both as a source of linguistic data (invaluable for investigating language-specific issues such as the diachronic evolution and dialectal differentiation of Greek and theoretical issues such as the mechanisms of language change and language contact) and as a depository of the cultural heritage of the now lost communities of the Greek-speaking populations of Asia Minor needs no comment. We will only say that Cappadocian was the only major dialectal variety of Modern Greek (the other being Pontic, Tsakonian, and Griko) which until now lacked a unified lexicographical treatment and that it is the only dialectal variety for which we now possess a wealth of new data, unavailable to previous scholarship, thanks to the recent re-discovery of 3rd generation native speakers. These two last factors were the original trigger for planning this project.

¹ See <http://cappadocian.upatras.gr/en> (access date 12.11.2022).

2. The electronic dictionary of the Cappadocian dialect

2.1. The dictionary and its sources

Turning to specific issues, we should firstly present the dictionary itself, its sources, format and content. The dictionary aspires to become the ultimate reference work on Cappadocian, and thus it attempts to incorporate and present all the available data on the Cappadocian language. This includes:

- (i) older written sources starting from the 19th century onwards (dictionaries, glossaries, linguistic descriptions, collections of primary texts such as folktales, songs, narrations, riddles etc.), most of which were collected and digitized thanks to an earlier project of the Laboratory of Modern Greek Dialects, AMiGre, and are available online².
- (ii) new oral recordings from current 3rd generation native speakers (descendants of Cappadocian refugees) collected the last decade. Special emphasis was placed on the exploitation of this new material, so that the dictionary under preparation is not a simple “super-container and comparative presenter” of already available but scattered data, but an opportunity for the presentation of new information. This allows also for a “diachronic” examination of the evidence, as we have the possibility to examine side-by-side data which may be divided by more than 100 years.

However, the dichotomy between written and oral sources has given rise to a major problem: whereas older sources are roughly equally distributed with respect to geographical provenance (i.e., data is available for almost all Cappadocian settlements, ca. 20 in number), oral data, from current speakers, are available only from 2–3 major communities, and mostly from that of Misti, which was the largest. This creates an imbalance in the lexicographical treatment of words, phenomena and senses. Another issue requiring special attention is the fact that the older material was in part collected by amateurs, or at a time when linguistic descriptive tools had not yet been sufficiently developed, and therefore it is to a certain extent

² See <http://amigredb.philology.upatras.gr/> (access date 12.11.2022).

less reliable than the oral material, containing many inaccuracies which can no longer be assessed.

Considering the large size of the corpus of data on which the dictionary is based, it is not surprising that the size of the dictionary is equally large: we are currently in the middle of letter S, and our lemma-list contains ca. 7000 entries. It is therefore projected that at the end of the project the dictionary will include ca. 8500 entries. This makes it much larger than the, similar in conception, dictionary of South Italian Greek (ca. 6.000 entries) but still smaller than the Tsakonian (ca. 13.000 entries) and the Pontic dictionaries (20.000 entries) — of course the data in each Cappadocian entry is much larger than in either of these last two and is similar to Karanastasis (1984). Luckily, “size does NOT matter” in the case of electronic publications, so we are letting ourselves run free in this respect. However, apart from the on-line edition, the “Historical Dictionary of the Cappadocian dialect” will also appear in print form, which will probably be realized as a four-volume-set.

The entry compilers of the Historical Dictionary of the Cappadocian dialect are, in alphabetical order, Gogo Katsouda, Io Manolessou, Symeon Tsolakidis, Christos Papanagiotou and George Chairetakis. Io Manolessou also acts as editor-in-chief, while Athanasios Karasimos and Io Manolessou are responsible for the creation of the dictionary’s DTD, i.e. the parametrization of the lexicographical software and the implementation of the dictionary’s electronic platform.

2.2. The e-dictionary

The format of the Dicadland dictionary attempts to conform to standards of state-of-the-art academic-level Dictionary Writing Systems (DWS), after careful evaluation of available options, and adheres to the most recent advances in the domain of electronic lexicography (see e.g. Granger & Paquot 2012). It is built using the powerful professional dictionary editing software TLex Suite, one of the most widely-used state-of-the-art DWS internationally. The decision to use a professional lexicographical tool, rather than one of the freely available online lexicographic platforms such as “Lexonomy” or “Matapuna” was forced by the size and complexity of the data to be treated, since free platforms do not leave much room for parametrization. On the other hand, the solution of creating from scratch an in-house custom-made lexicographic tool capable of handling our complex dialectal material was deemed impractical given the relative short time-range of the project.

In effect, the creation of a digital platform for the “Historical Dictionary of the Cappadocian Dialect” was something half- way between a retro-digitization and a new digitization proper (hence the title of the present paper). This is because we took the conscious decision to adhere the microstructure of the other major dialectal dictionaries of Greek we have mentioned above (Pontic, Tsakonian, South Italian), which in turn are all based on the “mother of all dialectal dictionaries”, the *Historical Dictionary of Modern Greek* (ILNE) of the Academy of Athens. Note that also the Dictionary of Medieval Greek of Kriaras (1968–) shares the same structure, as it too, was based, as a concept, on the *Historical Dictionary of Modern Greek*.

διακλύζω Σιν. **δακλύζω** Κερ.Κοτ.Τραπ.
Χαλδ. **δακλύω** Χαλδ. **δακλυῶ** Οίν. **δακλῶ**
Χαλδ. **δακλῶ** Χαλδ. **δακλύγω** Σάντ.Χαλδ.
Παθητ. **δακλύσκουμαι** Μετοχ. **δακλυσμένος**,
δακλυγμένος.
'Από τὸ ἀρχ. **διακλύζω**=ἐκπλύνω, ξεπλύνω.
1) Πλύνω καὶ ξεβγάλω τι μὲ καθαρὸ
νερό: **Δακλύζω** τὸ ποτήριον. 2) Καθαρίζω
σκεῦος τι διακινῶν αὐτὸ καὶ ἐκτινάσσω τὸ
περιεχόμενον ὕδωρ. 3) Πλύνω: **Δακλύζω** τὸ
στόμα μ'.

Picture 1: Papadopoulos (1958: 263)

***διάβασμα** το, **δέβασμα** ΒΧ, **δᾶβασμα** ΒΧ,
ζβᾶισμα ΜΠΤ, **ζιβᾶισμα** Μ.
Οἱ δύο τελευταῖοι τύποι ἀπὸ τον τύπ. **ζβαίχου**
<**διαβάζω**.
1. **Διάβασμα** κοιν. **Τὸ δέβασμα** σ' καλέ ρ,
δεβᾶζ' καλὰ Χ, **τὸ ζβᾶισμά** σι καλέρ ἐνί, ἐνί
ζβαίχου κα Μ, **τὸ διάβασμά** του καλό εἶναι,
διαβάζει καλὰ. 'Οσοί **έχου ζβαίγματα** σήμερε
Μ, **δεν έχεις διαβάσματα** σήμερα; 2. **Ευχή**
θρησκευτική κοιν. 'Ενί **θέου ζβᾶισμα**, **να μό-**
λει ο παπά να νι ζβαίσει Μ, **θέλει διάβασμα**
(ευχής), **νά 'ρθει ο παπάς να τον διαβάσει**.
3. **Καθοδήγηση**, **υποβολή συμπεριφοράς** κοιν.
Ἄν' εφωνιάε για ζβᾶισμα, **ὁ καταβήτσερε** Μ,
τον φώναξε να τον «διαβάσει», **δεν κατάλαβες;**

Picture 2: Costakis (1983: 262)

διάρμισμα τό, διάρμισμα [dǵármizma] Αιολ. (Κυδων.) Θράκ. ἑ. Ἑλλ. (Αἴν.) Ἰων. (Βουρλ.) Κάρπ. Κρήτ. Νάξ. (Ἀπίραθ.) Σέρρ. — Λεξ. Πρω. Δημητρ. Σταμ. Τσιούν. Ὑπερλ. διάρμ'σμα [dǵármizma] Ἰμβρ. Λέσβ. γάρμισμα [ǵármizma] Κάρπ. τζάρμισμα [dzármizma] Κάρπ. γί-άρμισμα [fǵármizma] Κάλυμν. διάρμοσμα [dǵármiozma] Κάρπ.

Ἀπό τὸ ρ. διαρμίζω (θ. ἄορ. διαρμισ-), ὅπου καὶ τύπ. γιαρμίζω, τζιαρμίζω, καὶ τὸ παραγωγ. ἐπίθμ. -μα. Ὁ τύπ. διάρμοσμα μὲ παρετυμολογ. ἐπίθρ. τοῦ ρ. ἀρμόζω.

1) Ἡ ἐνέργεια καὶ τὸ ἀποτέλεσμα τοῦ ρ. διαρμίζω 2, ἡ τακτοποίηση, ἡ τάξη Αἰολ. (Κυδων.) Θράκ. ἑ. Ἑλλ. (Αἴν.) Ἰμβρ. Ἰων. (Βουρλ.) Κάλυμν. Κάρπ. Κρήτ. Λέσβ. Σέρρ. — Λεξ. Πρω. Δημητρ. Σταμ. Τσιούν. Ὑπερλ.: *Θέλω καλὸ διάρμισμα στὸ σπίτι, ὄχι τσαπατσουλικά!* Ρεθύμν. (Μαλάκ.) *Μωρέ, χαρὰ στὸ γάρμισμα ἀπὸ 'καμὲς τοῦ σπιτίου!* Κάρπ. Συνών. βλ. λ. διάρμιση. Ἀντίθ. ἀδιαρμισιά. 2) Ἡ ἐνέργεια καὶ τὸ ἀποτέλεσμα τοῦ ρ. διαρμίζω 1, ἡ φύλαξη Ἰων. (Βουρλ.) 3) Ἡ ἐνέργεια καὶ τὸ ἀποτέλεσμα τοῦ ρ. διαρμίζω 3, τὸ ἀνακάτωμα Νάξ. (Ἀπίραθ.): *Ἦκαμὲς του πάλι σήμερα διάρμισμα τοῦ νεροῦ! Εἶδα λοῦται καὶ τὸ μετατόπισες ἂ' τὸ 'να μπιθάρει στ' ἄλλο;*

Picture 4b: From ILNE (2016: entry διάρμισμα)

δέμα(ν) το, Πιπράφρ. Μανασσ. Β 300, Ερμων. Κ 195, Δεζγ. πικδ. (Tsiouni) 627, Λιβ. Ρ 1427, Λιβ. Sc. 377, 752, Λιβ. Esc. 1494, 1844 (ζριτ. υπ.), Λιβ. Ν 1343, 1645, Αρχιλ. Ν 377, Ψευδο-Γεωργηλ., Ἄλ. Κων/π. 590, Σκλέντζα, Ποιήμ. 3^ο, Πεντ. Αρ. XXX 3, 13, 15, Στάθ. Γ' 439, Ροδολ. (Μανώσ.) Γ' [428], Ροδολ. Ε' [81], Συνομμ., Πιστ. φιδ. Γ' [831], Ε' [216], Τζάνε, Κρ. πάλ. 272²⁶. *δ ή μ μ ν*, Μιχλ. 20¹¹⁻¹³, 22¹², 252¹⁴⁻¹⁵, 654²⁸.

Το μτρν. ουσ. *δέμα*. Η λ. και σήμ. (Δημητράκ., λ. *δέμα*) και ως τοπων. (Βλ. Σαφρή I., Αθ. 40, 1928, 136).

1) Αυτὸ με τὸ ὁποῖο δένομε κ., σκεινί, τεινίχ (Η σημασ. μτρν., 1-5, λ. *δέμα* 1): *ἔλυσε με τὰ χέρια-της τοῦ πιττακίον το δέμαν* Λιβ. Sc. 377. 2) *Δέμα*, δεμάτι (Η σημασ. τον 5. αι., Lampe, Lex., λ. *δέμα* 1 και σήμ., Δημητράκ., ὁ.π. 3): *λάβρες, βεργόνια, δέματα καλὰ ἐταί σοβεμένα* Ροδολ. (Μανώσ.) Γ' [428]. 3) Δεσμός: *εσὺ το ἀδετον το δέμα της Τριάδος εἶσαι και μοροσόσιος αγάτη της μοιάδος*; Σκλέντζα, Ποιήμ. 3^ο: *μῆσα εἰς τη φιλιὰ ποῦ 'ναι σ' εσάς η τόση ἀτόμενε παντοτινὸ δέμαν ἀνάμεσάσας* Ροδολ. Ε' [81]. 4) Δεσμά: *'χ τὰ χέρια λῶσειτέ-την κι' ἀπὸ τ' ἀνάξια δέματα λυ-*

Picture 5: Kriaras (1977: 9)

The entry structure or microstructure of all the above-mentioned projects is, as you can see, tripartite: It consists of the following “sections”, which, were translated in the XML format of the Cappadocian dictionary’s template as “Elements” (for more information, see Karasimos et al. 2020):

- (i) A **Formal Section**, where the variant dialectal forms are set out, with phonetic transcription, part of speech characterization and geographical distribution, realized as “attributes”, and with bibliographic source data available as a drop-down list.
- (ii) An **Etymological section**, where the word’s origin (native/loanword) and dating are recorded. This is our only free-text field, as etymology remains a domain where XML encoding has yet to be standardized.
- (iii) A **Senses section**, with numbered senses and sub-senses, which include, as attributes, definitions, examples, quotations and documentation from oral and written corpora, and also including “special” types of examples, such as proverbs, songs and riddles. The latter type of examples required additional customization of the DTD, as they involved on the one hand “double” translations (literal and metaphorical) and on the other special formatting (verses).

3. Special issues of morphophonetic representation

3.1. Transcription

We have already touched upon the issue earlier, when we mentioned the problem that the older written sources of the dictionary employ a variety of symbolisms, not consistent with each other and not always easily interpretable. All of them had to a) be unified and homogenized and b) be represented in a system easily accessible and comprehensible. That is, the transcription needed to be given both in the International Phonetic Alphabet (for reasons of scientific clarity and accessibility to the international academic community) and in some form of the Greek alphabet (for reasons of accessibility to the general public).

In this respect, the project was greatly assisted by the previous work done on the topic by the *Historical Dictionary of Modern Greek* and its transcription system, which the Cappadocian historical dictionary also adopts. The specialized groundwork of the Historical Dictionary provides correspondence tables for all the variant transcription systems to be found in the standard dialectological publications on Cappadocian.

Καππαδοκική

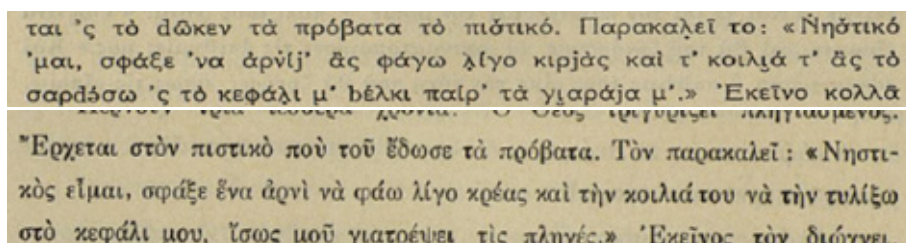
Καρακίτης 1885	Λαγαράς 1886	Πουκίτης 1906	Ασφραχίου 1946	Κετοργίου 1951	Μαυρογλυφίδης - Κίτσις 1960	Κωστήρας 1964	Αναστασιάδης 1976	Janse 2004, 2011	IANE	ΣΦΑ
		ä	ä		ä	ä	ä		ä	ae
		ö			ö		ö	ö	ö	e
		ü			ü	ü	ü	ü	ü	y
		ə		ə	ə	ə	ə	ɪ	i	u
			ä, ɪ, k̄	ä, ɪ, k̄	ä, ɪ, k̄	ä, ɪ, k̄	ä, ɪ, k̄		ä, ɪ, k̄	p', t', k'
			k̄	k̄	k̄	k̄	k̄		k̄	k
				k̄	k̄		k̄	c	k̄	c
		q							q	q
b, d, g		b, d, g	b, d, g	b, d, g	b, d, g	b, d, g	b, d, g	b, d, g	b, d, g	b, d, g
						ɣk̄	g		ɣk̄	g
	kh	ḡ	ḡ	ḡ	ḡ	ḡ	ḡ	x	ḡ	x
	ç		ḡ	ḡ	ḡ		ç		ḡ	c
	gh	ḡ	ḡ	ḡ	ḡ	ḡ	ḡ		ḡ	ɣ
	ç'		j	j	j	j	j	j		j
	λ'		λ	λ	λ	λ	λ		λ	l
			ú	ú	ú	ɣ, ú	ú	u	ú	n
			ú	ú	ú		ú	u	ú	n
		ɣ	ɣ (?)	ɣ (?)	ɣ (?)		ɣ (?)			ɣ
sch	oo	ö, ç	ö	ö	ö, ö	ö, ö	ö	ɪ	ö	f
		ç	ç	ç	ç	ç	ç	z	ç	s
ç'	no	no							no	ts
c'	τç	ç	τö	τö	τö, τö	τö	τö	ts	τö	ç'
g'	çç	j	τç	τç		τç	τç	dʒ	τç	dʒ

Picture 6: The variant transcription of Cappadocian texts and manuscripts (Manolessou et al. 2012)

If one is not aware of the potential variety of representation practices in the written sources, one can easily misread the data. To give a couple of concrete examples.

3.1.1. The case of Kentro Mikrasiatikon Spoudon approach

The system employed by the publications of the Kentro Mikrasiatikon Spoudon employ the symbols λ with superimposed or undersided dot, and ν with superimposed dot (Mavroxalybidis & Kesisoglou 1960).



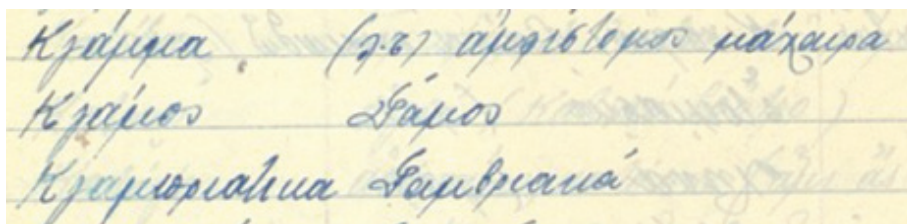
Picture 7: Text sample from Mavroxalybidis & Kesisoglou (1960)

This representation is in fact a trap: one is led to believe that it records a special kind of pronunciation of /l/ and /n/, when in fact it is nothing of the kind: these symbols represent the “normal” sounds [l] and [n]. The special notation starts from the assumption that the “default” realization of the phonemes /l/ and /n/ in Standard Modern Greek, when they are followed by the front vowels /e/ and /i/, is as the palatal allophones [ɫ] and [ɲ] — i.e. that Standard Modern Greek is like the local dialect of Patras or Eastern Crete. So a “special notation”, the dot in this case, is used to signal the ABSENCE of palatalization from Cappadocian. The dot just means that the words in question are supposed to be read: [paraka'li] [nifti'ko] [ar'ni] ['liyo] and not [paraka'li] [nifti'ko] [ar'ni] ['liyo].

3.1.2. Three Cappadocian sub-dialects

The second example is from three Cappadocian sub-dialects, those of Malakopi, Flogita and Silata, possess a realization of /k/ or /ɣ/ which is very velar, in fact [q]. This sound is represented in a variety of ways in the sources:

As <κ̣> with a superimposed dot, as plain <κ> or as <κγ>:



Picture 8: Text sample from Anthopoulos (1937–1938)

Κρέμασεν ένα κλωστή, και πήρεν ένα μήλο. Τόν δο έφαγεν, πόμη-
σο ρουργούρι τ. 'Απού άργάς τόν ήρταν τα έφτά παληκίρμα
πάλα το θύρα ήβραν δο ραπαδιμένο. "Τστερα άνοιξάν δο, κα.

Picture 9: Text sample from Dawkins (1916) (Silata)

All three representations <κ>, <κ̣> and <κγ>, are in fact the same sound, [q], which may correspond to two different phonemes, /k/ or /ɣ/. The *Historical Dictionary* of the Academy of Athens also offers a custom-made system of phonetic representation, based on the Greek alphabet, with some additional symbols used as diacritics, which combines on the one hand systematic one-to-one mapping with the IPA and at the other increased readability and comprehensibility through the retention of standard Greek orthography. We cannot stress how crucial this is for dialects such as Cappadocian, where massive phonetic and morphological changes have rendered even common words of the standard completely unrecognizable.

For example, it is by far preferable to spell *αυιτσα* or *ευιτσα* instead of *αβιτσα* /avitsa, since the etymological connection to *αυγή*, *αυγίτσα* becomes immediately obvious and the sense ‘dawn’, ‘morning’ transparent. Similarly, *νεγγριώνα* spelt with double <γγ> instead of <γκ> or <g> provides an insight to the origin of the word, that it has something to do with *αγγούρια* ‘cucumbers’, and thus can help one to guess its meaning ‘vegetable garden’

Of course, the use of diacritics adds a level of difficulty in digitization, as these are not always easily rendered with the standard fonts available, especially when they stack one on top of the other. Whereas the *Historical Dictionary of Modern Greek* employs a custom-made font of its own, **Athens Academy fonts**, it was a conscious decision early on in our planning not to use any fonts or glyphs which are not part of the Unicode Standard, and to employ only fonts which belong to Microsoft’s ClearType collection (such as Cambria, Calibri, Corbel) because these allow for much better alignment of diacritics on top of main letter glyphs.

4. Lemmatization

One of the most serious problems we face in compiling the Cappadocian dictionary is that of lemmatization, i.e. the selection of the headword. This is because this dialect has undergone a great number of phonetic and morphological changes, more than any

other Greek dialectal variety, both due to the length of time of its separation / isolation from the other forms of Greek, about a thousand years, (from the 11th c.), and due to intense contact with Turkish. This entails that a) its forms frequently are very distanced from the Standard and most other dialectal varieties of Greek and b) more importantly, each local Cappadocian form is distanced from all others; each settlement may present its own variant (see also Manolesou & Katsouda 2021; Forthcoming).

Let's give a concrete example:

καλλικεύω (ρ.) *καλγκεύω* [kaɫˈʝeɒ] *Αξ. καλγεύω* [kaɫˈʝeɒ] *Σίλατ. γαλγεύω* [ɣaɫˈʝeɒ] *Σινασσ. -Αλεκτορίδης 1883/5 καλιτσέου* [kaɫˈʝeɒ] *-Κοτσανίδης 2006 γκαλιτσέου* [ɣaɫˈʝeɒ] *Φάρασ. γαλιτσέου* [ɣaɫˈʝeɒ] *Φάρασ. -Παπαστεφάνου/Καρακαλιδίου 2009 καλιεύω* [kaɫˈleɒ] *Φλογ. καλιεύω* [kaɫˈleɒ] *Μαλακ., Μισθ. καλιεύω* [kaɫˈleɒ] *Μισθ. -Dawkins 1916 γαλιεύω* [ɣaɫˈleɒ] *Σινασσ. -Αρχιλαός 1899 γαλιεύω* [ɣaɫˈleɒ] *Σίλατ. καλντεύω* [kaɫˈdeɒ] *Αξ. Αραβαν., Γούρδ., Τελμ., Φερτάκ. γκαλντεύω* [ɣaɫˈdeɒ] *Ουλαγ. κατλεύω* [kaɫˈleɒ] *Σιλ. Αόρ. καλλικεφα* [kaɫˈlikeɓa] *Σίλατ. -Χωλόπουλος 1897 καλλιτζεφα* [kaɫˈliɟeɓa] *Φάρασ. γκαλλιτζεφα* [ɣaɫˈliɟeɓa] *Φάρασ. -Dawkins 1916 γαλλιτζεφα* [ɣaɫˈliɟeɓa] *Φάρασ. -Παπαδόπουλος 2012 κάλγκεφα* [kaɫˈgeɓa] *Αξ. κάλγεφα* [kaɫˈgeɓa] *Ποτάμ., Σίλατ. κάλντεφα* [kaɫˈdeɓa] *Τελμ.*
Από το μεσν. καλλικεύω τύπ. του ρ. καβαλλικεύω, πβ. Πικατ. 365 "να καλλικεύω το φαρίν, τον κόσμο να γυρίσω". Πβ. και Λεξ. Σομαβέρα, λ. καβαλλικεύω, όπου και τύπ. καλιεύω. Πιθ. οι τύπ. καταγραφόμενοι ως καλγ-, γαλγ- να αποδίδουν πραγμάτωση [kaɫg], [ɣaɫg], οπότε η εμφανής εξέλιξη είναι καλλικεύω > καλ'κεύω > καλγεύω με αποβολή του άτονου [i] και ηχηροποίηση του [k] > [g] μεταξύ ηχηρών φθόγγων. Αυτό μάλλον υπονοεί ο Dawkins (1916: 605), λέγοντας ότι ο πλησιέστερος στο αρχ. έτυμο είναι ο τυπ. αορ. κάλγεφα. Η εξέλιξη καβαλλικεύω > καλλικεύω σχολιάζεται στο Ε. Κριαρά, Η Ρίμα θρηνητική του Ιωάννου Πικατόρου, *Επετηρίς του Μεσαιων. Αρχείου* 2 (1940), Γλωσσάριον, σ. 66. Ο τυπ. καλλικεύω και σε νοτιοανατολικά ιδιώμ. (Κάρπ., Κόσ., Χίος)

Picture 10: The entry καλλικεύω from TLex software entry

For the verb ‘to ride’ each settlement presents different variants, such as *καλγκεύω*, *καλτσέου*, *καλιεύω*, *γαλιεύω*, *γκαλντεύω*, *κατλεύω*. The only way to unify them is to subsume them under an unattested headword *καλλικεύω*, which corresponds to the word’s medieval etymon, the point of origin of all the changes (n.b. in turn this comes from the Standard word *καβαλλικεύω*). Whether the headword will be an unattested reconstructed one or not to a certain extent depends on the available evidence of the sources. One of the Cappadocian dialects, that of Sinasos, is much more conservative than all the others (or much more affected by the standard?). If the word happens to be preserved in the texts from Sinasos, then the headword will probably be attested. For example:

αγκάθι (ουσ. ουδ.) *αγκάθι* [aŋˈgaθi] *Σινασσ. αγκάθ'* [aŋˈgaθ'] *Μαλακ., Σίλατ., Φλογ. αγκάγ'* [aŋˈgaɣ'] *Αξ. αγκάχ'* [aŋˈgaχ'] *Μισθ. αγκάζ'* [aŋˈgaɟ'] *Σεμέντρ. αγκάτ'* [aŋˈgaɪ'] *Φερτάκ. αγκάρ'* [aŋˈgaɾ'] *Αραβαν., Γούρδ. γκάθι* [ŋgaθi] *Φάρασ. Μεσν. ουσ. άγκάθι το οπ. από το αρχ. άκίθιον.*

Picture 11: The entry αγκάθι and its variant formations

The problem becomes more acute when the phonetic change involves the initial vowel or consonant of the word, as this affects the alphabetical order of the dictionary (this is less of a problem in an electronic than in a print dictionary, but still there are many readers who “browse” instead of just “search”).

A concrete example involves the adaptation of loanwords from Turkish which start with /k/ (e.g. *καμπούρης* ‘hunchback’ from Turk. *kambur*). Since the local Turkish dialects of Eastern Anatolia present a change /k/ > /ɣ/, loanwords entering Cappadocian from local Turkish varieties often show initial /ɣ/, e.g. *γαμπούρ*’ Misti.

καμπούρης (επίθ.) *καμπούρης* [kam'buris] Γούρδ. *γαμπούρ*' [ɣam'bur] Μαλακ. *γαμπούρ* [ɣam'bur] Μισθ. Φάρασ. *γαμπύρ* [gam'bur] Ουλαγ. θηλ. *γαμπουρούτσα* [ɣambu'rutʃa] Φάρασ.
Από το τουρκ. επίθ. *kambur* = *καμπούρης*.

Picture 12: The entry *καμπούρης* και its variant formations with initial /k/ and /ɣ/

If there exists at least one subdialect which does not present this phenomenon, the case is straightforward, the headword will be an attested word starting with <κ>. But if not, then the headword, unless one wishes to resort to reconstruction again (less justifiable in the case of loanwords), will start with <ɣ>.

γαπαχλούς (ουσ. αρσ.) *γαπαχλούς* [ɣapax'lus] Φάρασ.
Από το τουρκ. επίθ. *kabaklı*.
1 Κολοκύθι ως γλυκό κουταλιού
γαπαχουφαϊ (ουσ. ουδ.) *γαπαχουφαϊ* [ɣapaxu'fa'i] Φάρασ.
Από τα ουσ. *καμπάκι*, όπου και τύπ. *γαπάχι*, και *φαϊ*.
Φαγητό με κολοκύθια.

Picture 13: The entries *γαπαχλούς* and *γαπαχουφαϊ*

Here the derivative and the compound are to be found only in Pharasa, and so start with <ɣ>. However, the simplex word is to be found in many areas, and therefore the headword starts with <κ>, something which disrupts the word-family.

καμπάκι *γαμπάκι*' [gabak] Ουλαγ. *γαμπάχ'* [ɣa'bak] Μαλακ., Ποτάμ. *γαμπάχ'* [ɣa'bak] Αραβαν., Μισθ. *γαπάχι* [ɣa'paxi] Φάρασ. Πληθ. *γαμπάκια* [ɣa'baʃa] Μισθ.
Από το τουρκ. ουσ. *kabak* = κολοκύθι.
1 Κολοκύθι ό.π.τ.: *Γαπαχού γούτσια* (Σπόρια κολοκυθίου) Μισθ. -Κοτσανίδης 2006 *Εγώ μι γαπάια να σπέρου* (Εγώ με κολοκύθια θα σπείρω) Μισθ. -ΑΠΥ-Καρατο. *Έχου γαμπάχια, ε δάνου δα φέτις φέτις, τεγανίζου δα τρώου' μετά εκείνου του γαμπάχ' ξίνου δου, ξίνου δου, μάζου λίου αλεύρι' δάνου κιοφτάδης* (Έχω κολοκύθια, τα κάνω φρέτες-φέτες, τα τηγανίζω τα

Picture 14: The entry *καμπάκι*

It really is questionable whether one should create an unattested form such as “καμπακοφάι” in order to keep the family together. As in the previous case, the issue is again connected with the availability of attestations. If we had more sources discussing food in other Cappadocian villages, perhaps our lemmatization would be different, as some forms with initial /k/ may have been available.

5. Concluding remarks

In the short time slot available to us we have hoped to show that the Dicadland project is a major, state-of-the-art endeavour in the field of dialectal lexicography, which will produce a concrete and highly useful work of reference. The data discovered and presented allow for important advances in our knowledge of the Greek dialects. This includes both the interpretation of known data (homogenization, phonetic analysis, new and improved etymologies, explanation of changes, establishment of semantic ranges through comparative examination) and the addition of new and original data. The methodology adopted brings to the fore many thorny issues of dialectal lexicography, which we hope, if not to solve, at least to highlight from several aspects.

To conclude, we provide a golden standard template for future dialectal dictionaries, computational lexicography and digital humanities and we introduce a new era of Digital Humanities in Dialectology and Lexicography by following the dominant trends and schemata of this new discipline. Moreover, we build a concrete and fully adaptable basis for the preservation of endangered dialects of Modern Greek and their cultural heritage. Finally, the open data — open access policy will help the extension of the research community and engage the native speakers to provide more data in the near future.

References

- Anthopoulos, Th. 1937–1938. *Φλογητά* [*Flogita*]. Unpublished manuscript, ILNE archive No. 1547. Athens: Academy of Athens [in Greek].
- Costakis, Th. 1986. *Λεξικό της τσακωνικής διαλέκτου* [*Dictionary of the Tsakonian Dialect*]. Vol. 1. Athens: Academy of Athens [in Greek].

- Dawkins R. M. 1916. *Modern Greek in Asia Minor. A Study of the Dialects of Silli, Cappadocia and Phárasa with Grammar, Texts, Translations and Glossary*. Cambridge: Cambridge University Press.
- Granger, M. & S. Paquot. 2012. *Electronic Lexicography*. Oxford: Oxford University Press.
- ILNE 1933 — *Ιστορικών Λεξικόν τῆς Νέας Ἑλληνικῆς, τῆς τε κοινῶς ὁμιλουμένης καὶ τῶν ιδιωμάτων* [*Historical Dictionary of Modern Greek, both of the Standard and the Dialects*]. Vol. 1. Athens: Academy of Athens [in Greek].
- ILNE 2016 — *Ιστορικών Λεξικόν τῆς Νέας Ἑλληνικῆς, τῆς τε κοινῶς ὁμιλουμένης καὶ τῶν ιδιωμάτων* [*Historical Dictionary of Modern Greek, both of the Standard and the Dialects*]. Vol. 7. Athens: Academy of Athens [in Greek].
- Innsbruck EDD Online 3.0 — English Dialect Dictionary, University of Innsbruck. Available at: <http://eddonline-proj.uibk.ac.at/edd/termsOfUse.jsp> (accessed on 28.11.2022).
- Karanastasis, A. 1986. *Ιστορικόν λεξικόν των ελληνικῶν ιδιωμάτων της Κάτω Ιταλίας* [*Historical Dictionary of the Greek Dialects of South Italy*]. Vol. 2. Athens: Academy of Athens [in Greek].
- Karasimos, A., I. Manolessou I. & D. Melissaropoulou. 2020. Creating a DTD template for Greek dialectal lexicography: the case of the Historical Dictionary of the Cappadocian dialect. In Z. Gavriilidou, M. Mitsiaki & A. Fliatouras (eds), *Proceedings of the 19th Congress of the European Association for Lexicography- EURALEX XIX, Alexandroupolis, September 2021*. Vol. 1. Komotini: Democritus University of Thrace, 305–314.
- Kriaras, E. 1968-. *Λεξικό τῆς μεσαιωνικῆς Ἑλληνικῆς δημώδους γραμματείας, 1100–1669* [*Dictionary of Greek Medieval Vernacular Literature, 1100–1669*]. Vol. 1–22. Thessaloniki: Centre for the Greek Language [in Greek].
- Kriaras, E. 1977. *Λεξικό τῆς μεσαιωνικῆς Ἑλληνικῆς δημώδους γραμματείας, 1100–1669* [*Dictionary of Greek Medieval Vernacular Literature, 1100–1669*]. Vol. 5. Thessaloniki: Centre for the Greek Language [in Greek].
- Manolessou, I. & G. Katsouda. 2021. Drawing the line between synchrony and diachrony in historical and dialectal lexicography. In Z. Gavriilidou, L. Mitits & Sp. Kiosses (eds.), *Proceedings Book of Euralex XIX (Congress of the European Association for Lexicography)*. Vol. 2. Komotini: Democritus University of Thrace, 83–92.
- Manolessou, I. & G. Katsouda. Forthcoming. On lemmas and dilemmas again: problems in historical dialectal lexicography. In *Proceedings of the 11th International Conference on Historical Lexicography and Lexicology (Universidad de La Rioja, 15–18 June 2021)*.
- Manolessou, I., S. Beis & C. Basea-Bezantakou. 2012. *Η φωνητική απόδοση των νεοελληνικῶν διαλέκτων* [*The phonetic transcription of the modern Greek dialects*]. *Λεξικογραφικόν Δελτίον* [*Lexicographic Bulletin*] 26: 161–222 [in Greek].
- Mavrochalyvidis, G. & I. Kesisoglou. 1960. *Το γλωσσικό ιδίωμα της Αξού* [*The Dialect of Axos*]. Athens: Centre for Asia Minor Studies.
- Papadopoulos, A. A. 1958. *Ιστορικόν λεξικόν τῆς ποντικῆς διαλέκτου* [*Historical Lexicon of the Pontic Dialect*]. Vol. 1: Α–Λ. Athens: Epitropi Pontiakon Meleton [in Greek].