

# Προκλήσεις επισημείωσης ενός πολυ-διαλεκτικού, πολυ-επίπεδου σώματος γραπτών και προφορικών κειμένων των Νεοελληνικών Διαλέκτων

Αθανάσιος Καρασίμος<sup>1,3</sup>, Ελένη Γαλιώτου<sup>2</sup>, Νικήτας Καρανικόλας<sup>2</sup>, Γιώργος Κορωνάκης<sup>2</sup>, Κώστας Αθανασάκος<sup>2</sup>, Δημήτρης Παπαζαχαρίου<sup>1</sup> & Αγγελική Ράλλη<sup>1</sup>  
Πανεπιστήμιο Πατρών<sup>1</sup>, ΤΕΙ Αθηνών<sup>2</sup>, Ακαδημία Αθηνών<sup>3</sup>

## 1. Εισαγωγή

### 1.1 THALIS project AMiGre

Στην παρούσα μελέτη, η οποία αποτελεί μέρος του προγράμματος «AMIGRE-Πόντος, Καππαδοκία, Αϊβαλί: στα χνάρια της Μικρασιατικής Ελληνικής Γλώσσας», παρουσιάζεται η επισημείωση ενός διαλεκτικού σώματος αρχείων που διαφέρει από τα υπόλοιπα στα δύο ακόλουθα βασικά σημεία. Γίνεται συστηματοποιημένη προσπάθεια επισημείωσης με κοινή στρατηγική σε γραπτά και προφορικά δεδομένα που αφορούν σε ένα μεγάλο εύρος δειγμάτων από τις διαλεκτικές ποικιλίες του Πόντου, της Καππαδοκίας και του Αϊβαλιού. Πιο συγκεκριμένα, σκοπός του ερευνητικού προγράμματος είναι να μελετήσει συστηματικά τα Ποντιακά, τα Καππαδοκικά και τα Αϊβαλιώτικα, τρεις γλωσσικές ποικιλίες που απειλούνται με εξαφάνιση. Μεταξύ άλλων, επιδιώκεται η μελέτη των συγκεκριμένων διαλέκτων με σκοπό να αποκαλυφθούν οι ομοιότητες και οι διαφορές τους σε συγχρονικό επίπεδο, να επισημανθεί η εξέλιξή τους, να χαρτογραφηθεί η διαφοροποίησή τους, αλλά και να εντοπισθούν τα σημαντικότερα χαρακτηριστικά τους σε σχέση με τις υπόλοιπες Νεοελληνικές διαλέκτους. Επιπροσθέτως, γίνεται προσπάθεια για μία εμπειριστατωμένη ανάλυση συγκεκριμένων φωνητικών/φωνολογικών, μορφολογικών και σημασιολογικών φαινομένων, καθώς και της επιρροής διαφορετικά τυπολογικών γλωσσικών συστημάτων, μιας και είναι εμφανής η επίδραση της Τουρκικής (συγκολλητικής γλώσσας), στις συγκεκριμένες διαλέκτους της Νέας Ελληνικής (διαχτυτικής γλώσσας). Για αυτό το λόγο έχει γίνει συστηματική αρχειοθέτηση και ψηφιοποίηση προφορικού και γραπτού υλικού μεγάλου εύρους και έχει οργανωθεί σε μία ψηφιακή βάση δεδομένων. Ένα σημαντικό μέρος του πρωτογενούς υλικού θα μεταγραφεί και θα σχολιαστεί με την χρήση του πιο σύγχρονου εξοπλισμού. Γραπτό υλικό θα ψηφιοποιηθεί, και ένα μέρος αυτού, που θα επιλεγεί σύμφωνα με αυστηρά ποιοτικά κριτήρια (χρονολόγηση, προέλευση, αξιοπιστία), θα μεταγραφεί.

### 1.2 Σώματα γραπτών κειμένων vs Σώματα προφορικών κειμένων

Η συνέπεια στην επισημείωση σωμάτων κειμένων είναι μια ουσιώδης ιδιότητα για τις πολλαπλές χρήσεις επισημειωμένων σωμάτων κειμένων στην υπολογιστική και θεωρητική γλωσσολογία. Παλαιότερες έρευνες εντόπισαν προβλήματα σε μορφολογική και POS επισημείωση (van Halteren 2000· Eskin 2000· Dickinson & Meurers 2003), ενώ πιο πρόσφατες εντόπισαν λάθη σε συντακτικό και δομικό επίπεδο (Ule & Simon 2004· Dickinson 2005).

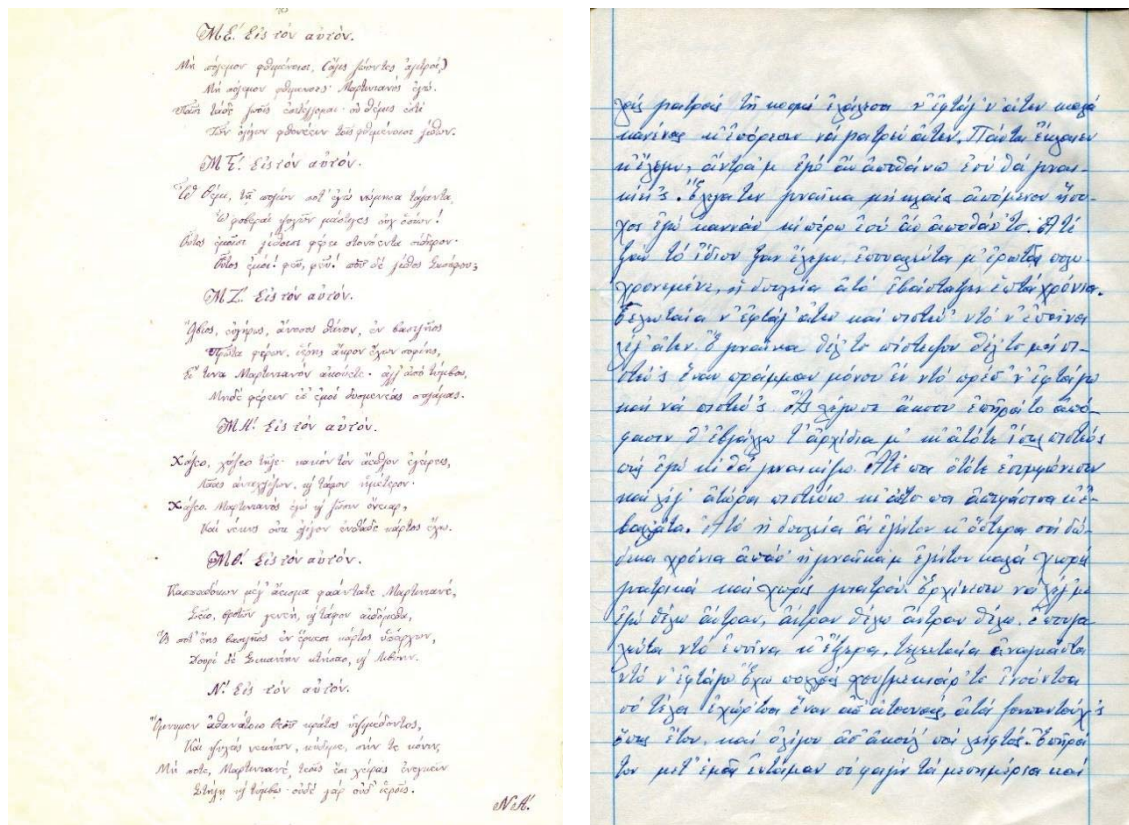
Τα σώματα γραπτών κειμένων είναι σαφώς περισσότερα ανά γλώσσα παγκοσμίως και σημαντικό ερευνητικό κομμάτι της υπολογιστικής και διακειμενικής γλωσσολογίας στοχεύει στην επισημείωση και την αξιοποίησή τους. Από την άλλη, τα σώματα προφορικών κειμένων υστερούν σε όγκο και διαφέρουν σε πολλά σημεία από τα αντίστοιχα γραπτά, ωστόσο υπάρχει έλλειψη συντονισμένης επισημείωσης, ενώ το ζήτημα της ανίχνευσης σφαλμάτων στον σχολιασμό της ομιλούμενης γλώσσας σωμάτων δεν έχει ακόμη αντιμετωπιστεί συστηματικά. Αυτό είναι σημαντικό δεδομένου ότι τα σώματα προφορικών κειμένων αυξάνονται ιδιαίτερα, όπως φαίνεται στο Linguistic Data

Consortium ([www ldc.upenn.edu](http://www ldc.upenn.edu)). Το πρόβλημα εντείνεται όταν γίνεται προσπάθεια δημιουργίας κοινής στρατηγικής επισημείωσης σε σώματα προφορικών και γραπτών κειμένων και δη όταν το αντικείμενο είναι ιδιαίτερα εξειδικευμένο, όπως το προαναφερθέν διαλεκτικό σώμα.

## 2. State-of-the-art: σχεδιασμός συστήματος

### 2.1 Η φύση των δεδομένων

Το σώμα προφορικών κειμένων του έργου AMiGre αποτελείται από ηχογραφήσεις περίπου 180 ωρών (δηλαδή 60 ώρες ανά διάλεκτο), όπως αυτές συλλέχθηκαν για τη διαλεκτική βάση Gree.D. (Karasimos et al. 2008). Η συλλογή των ηχογραφήσεων έγινε με συσκευές ψηφιακής ηχογράφησης υψηλής ευκρίνειας, σε όσον το δυνατόν πιο ήσυχες συνθήκες και πάντα με συναίνεση των συνομιλητών. Η επιλογή των ομιλητών έγινε με μεγάλη προσοχή, όσο αυτό ήταν εφικτό· στόχος ήταν οι ομιλητές να έχουν καθαρή άρθρωση, να έχουν φυσική ροή ομιλίας, και να κάνουν συστηματική χρήση της διαλέκτου στην καθημερινότητά τους. Επίσης, στις περισσότερες περιπτώσεις ήταν απαραίτητη η ύπαρξη του ενδιάμεσου στις ηχογραφήσεις, ώστε οι ομιλητές να αισθάνονται πιο οικεία κατά την διάρκεια της ηχογράφησης και να ελαχιστοποιηθούν τα σημεία διαλόγου όπου θα γινόταν αλλαγή γλωσσικού συστήματος επικοινωνίας (εγκατάλειψη της διαλέκτου και χρήση της Κοινής Νέας Ελληνικής). Βασική προϋπόθεση για τον ενδιάμεσο ήταν η καλή σχέση και γνωριμία με τους ομιλητές καθώς και η άριστη γνώση και χρήση της διαλέκτου.



Εικόνα 1 & 2: Δείγμα εικόνων από τα ψηφιοποιημένα χειρόγραφα (αριστερά Επιτάφια επιγράμματα του Λεβίδη· δεξιά Χειρόγραφα Καζαντζίδη).

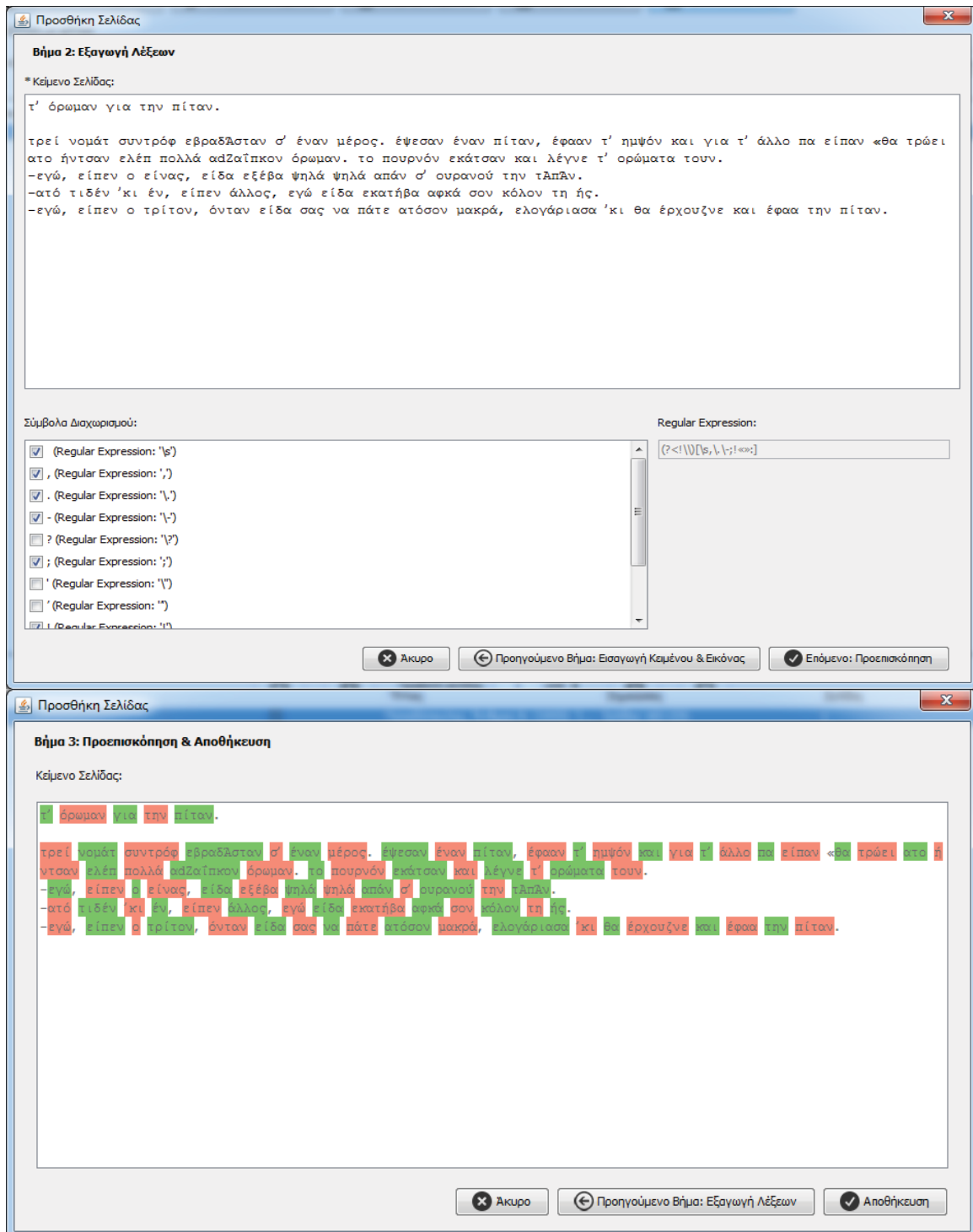
Αντιστοίχως, το σώμα γραπτών κειμένων αποτελείται από ψηφιοποιημένα χειρόγραφα συνόλου 2.000.000 λεξικών τύπων. Το σημαντικότερο ζήτημα για τη συλλογή γραπτών δεδομένων είναι η έλλειψη πρωτογενών πηγών και κυρίως χειρογράφων για τα Αϊβαλιώτικα· αναπόφευκτα η ισορροπία

ανάμεσα στην αντιπροσωπευτικότητα του δείγματος κειμένων που ψηφιοποιήθηκαν δεν ήταν εφικτή. Πέραν αυτής της εγγενούς δυσκολίας, τα κείμενα επιλέχθηκαν με βάση συγκεκριμένα κριτήρια. Βασικό κριτήριο ήταν το ζήτημα πνευματικής ιδιοκτησίας για την ψηφιοποίηση και για αυτό το λόγο επιλέχθηκαν κείμενα πριν το 1938. Επίσης, επιλέχθηκαν κυρίως κείμενα πεζού λόγου με ελάχιστη επιλογή ποιημάτων και τραγουδιών. Εκτός από ένα αντιπροσωπευτικό δείγμα ανάμεσα στα δημοσιευμένα κείμενα και τα χειρόγραφα, σημαντικό βάρος δόθηκε στη σπανιότητα μερικών εξ αυτών (αναλυτικά για τα κριτήρια στο Κολιοπούλου, Μαρκόπουλος & Παντελίδης 2015). Τα δεδομένα των παραπάνω σωμάτων πέρασαν από επεξεργασία, επιλογή, επισημείωση και ανάλυση και επεξεργάζονται σύμφωνα με το μοντέλο 3A (Annotation, Abstraction, Analysis) των Wallis & Nelson (2001) και τον προτεινόμενο μορφότυπο των Gries & Berez (υπό έκδοση). Για την περαιτέρω επεξεργασία, επισημείωση, ανάλυση και περιγραφή μεταδεδομένων έγιναν δύο υποσώματα κειμένων με 60 ώρες και 200.000 λέξεις αντίστοιχα. Η συγκεκριμένη επεξεργασία και ανάλυση έγινε εκτός από τη συνδρομή δημοφιλών γλωσσολογικών εργαλείων, με επτά νέες εφαρμογές που δημιουργήθηκαν στο πλαίσιο του προγράμματος (βλ. ενότητα 2.2).

## 2.2 Οι εφαρμογές του συστήματος

Το σύστημα διαθέτει επτά (7) βασικές εφαρμογές για την υποστήριξη της ανάλυσης των συγκεκριμένων διαλεκτικών σωμάτων, ενώ παράλληλα γίνεται η χρήση δύο εξαιρετικά δημοφιλών γλωσσολογικών εργαλείων, όπως είναι το Praat και το ELAN. Οι επτά (7) αυτές βασικές εφαρμογές είναι οι ακόλουθες (αναλυτικότερα βλ. Karanikolas, Galiotou & Ralli 2014):

- (α) **Phon Tagger** για την οριοθέτηση των λέξεων. Χρησιμοποιείται τόσο στο προφορικό όσο και στο γραπτό σώμα κειμένων, ώστε να υπάρχει μια ενιαία αντιμετώπιση της πληροφορίας των μορφολογικών ορίων των λέξεων μεταξύ των δύο σωμάτων.
- (β) **Morph Tagger** για τον μορφολογικό σχολιασμό των λέξεων, όπου πραγματοποιείται στο επίπεδο λέξης. Για κάθε μορφολογική λέξη παρέχονται πληροφορίες σχετικά με το μέρος του λόγου, γραμματικές ιδιότητες και μορφολογικά φαινόμενα, όπως η παραγωγή και η σύνθεση.
- (γ) **Synt Tagger** για τη συντακτική ανάλυση και δομή φράσεων και προτάσεων· στην τρέχουσα κατάσταση του συστήματος, η επισημείωση γίνεται σε επίπεδο λέξης, όπου κάθε λέξη συνδέεται τουλάχιστον με μία συντακτική δομή. Η εφαρμογή παρέχει επίσης η δυνατότητα για επισημείωση σε μια φράση ή σε προτασιακό επίπεδο.
- (δ) **Sem Tagger** για το σημασιολογικό σχολιασμό. Καταχωρούνται πληροφορίες όπως *δάνειο* (καθώς και την καταγωγή του), *ιδιωματική φράση*, κτλ.
- (ε) **Text Imaging** για την προεπισκόπηση εικόνων από τα ψηφιοποιημένα κείμενα και χειρόγραφα.
- (στ) **Text Transcription** για μεταγραφή των ψηφιοποιημένων κειμένων και των εικόνων.
- (ζ) **MOS (Oral Metadata)** για μια ολοκληρωμένη δομή μεταδεδομένων· αυτή η εφαρμογή παρέχει τη δυνατότητα διατήρησης και ενημέρωσης των μεταδεδομένων του σώματος προφορικών κειμένων και περιλαμβάνει πληροφορίες όπως *ηλικία*, *φύλο*, *πολιτισμικό υπόβαθρο* του ομιλητή μεταξύ άλλων (σημειώνεται ότι, οι πληροφορίες αυτές δεν είναι διαθέσιμες για τις γραπτές πηγές).



Εικόνα 3 & 4: Δείγμα από την εφαρμογή-υποσύστημα οριοθέτησης μορφολογικών λέξεων.

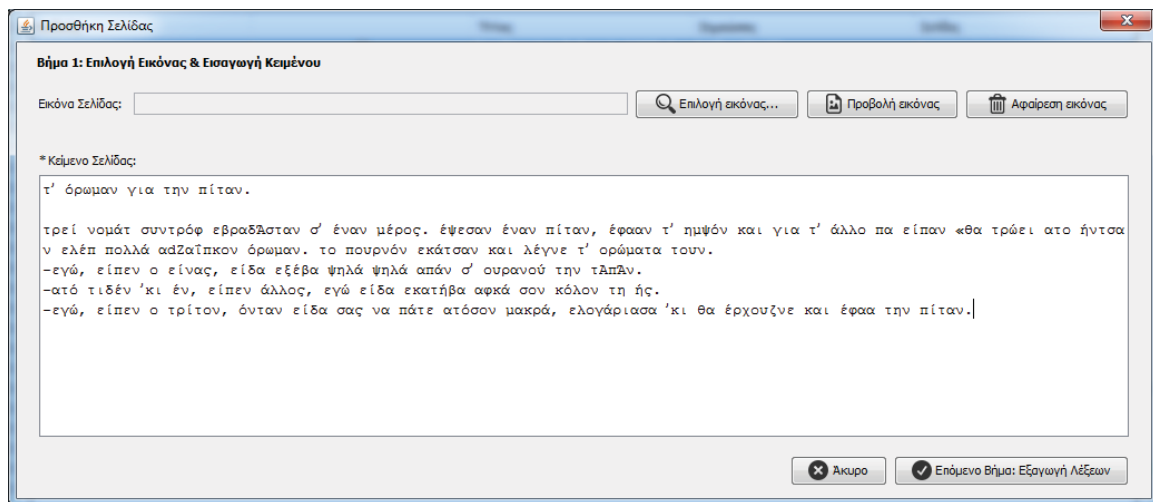
### 3. Προ-επεξεργασία σωματών γραπτών και προφορικών κειμένων

Η προ-επεξεργασία των δεδομένων μπορεί να συνοψιστεί ως εξής:

- (α) **Διαμόρφωση και Παραμετροποίηση:** Το κάθε πρωτογενές στερεοφωνικό ηχητικό αρχείο διαχωρίστηκε στα αντίστοιχα κανάλια του και έγινε επιλογή των κατάλληλων αρχείων με βάση συγκεκριμένα γλωσσολογικά και τεχνικά κριτήρια (βλ. Karasimos et al. 2008).

Επιπροσθέτως, οι εικόνες πέρασαν από τεχνική επεξεργασία για απομόνωση των σελίδων, αποκοπή μαύρων πλαισίων και ρύθμιση της καθαρότητάς τους.

- (β) **Επισημείωση:** Το σώμα γραπτών κειμένων πέρασε από μια συστηματική παραμετροποίηση φωνολογική και μορφολογική με βάση προεπιλεγμένες ετικέτες για ελεγχόμενες λίστες τιμών για την πλήρη κάλυψη των δύο επιπέδων. Παράλληλα κωδικοποιήθηκε μια μικρή παραλλαγή του προτύπου SAMPA (Wells 1997) και ενοποιήθηκαν οι διαφορετικές ποικιλίες συμβόλων γραπτών κειμένων με βάση τη πρόταση των Μανωλέσσου, Μπέη και Μπασέα (2012). Για τη επεξεργασία των προφορικών κειμένων έγινε μια αρχική προετοιμασία σύμφωνα με μια ανανεωμένη προσέγγιση παλαιότερης τακτικής επισημείωσης (Ράλλη, Παπαζαχαρίου & Καρασίμος 2010). Συγκεκριμένα, από το σύνολο του ψηφιοποιημένου υλικού οι επιλεγμένες λέξεις μεταγράφηκαν «δια χειρός», χωρίς την βοήθεια αυτοματοποιημένου λογισμικού μεταγραφής, λόγω των δυσκολιών που ένα τέτοιο εγχείρημα ενδεχομένως να προκαλούσε, όπως είναι η δυσκολία αυτόματης αναγνώρισης πολυτονικού συστήματος, δυσκολία αυτόματης αναγνώρισης χαρακτήρων στο χειρόγραφο υλικό (Κολιοπούλου, Μαρκόπουλος & Παντελίδης 2015).



Εικόνα 5: Δείγμα από την επισημείωση κειμένου ενός χειρόγραφου.

- (γ) **Μεταδεδομένα:** Ακολουθήθηκε το πρωτόκολλο καταγραφής για τα προφορικά δεδομένα, όπου επιλέχθηκαν οι πληροφορίες που ταιριάζουν και για τα σώματα γραπτών κειμένων με την παράλληλη εισαγωγή νέων ελεγχόμενων λιστών με τιμές για τα ψηφιοποιημένα κείμενα.

## 4. Επισημείωση

### 4.1 Επισημείωση σώματος γραπτών και προφορικών κειμένων

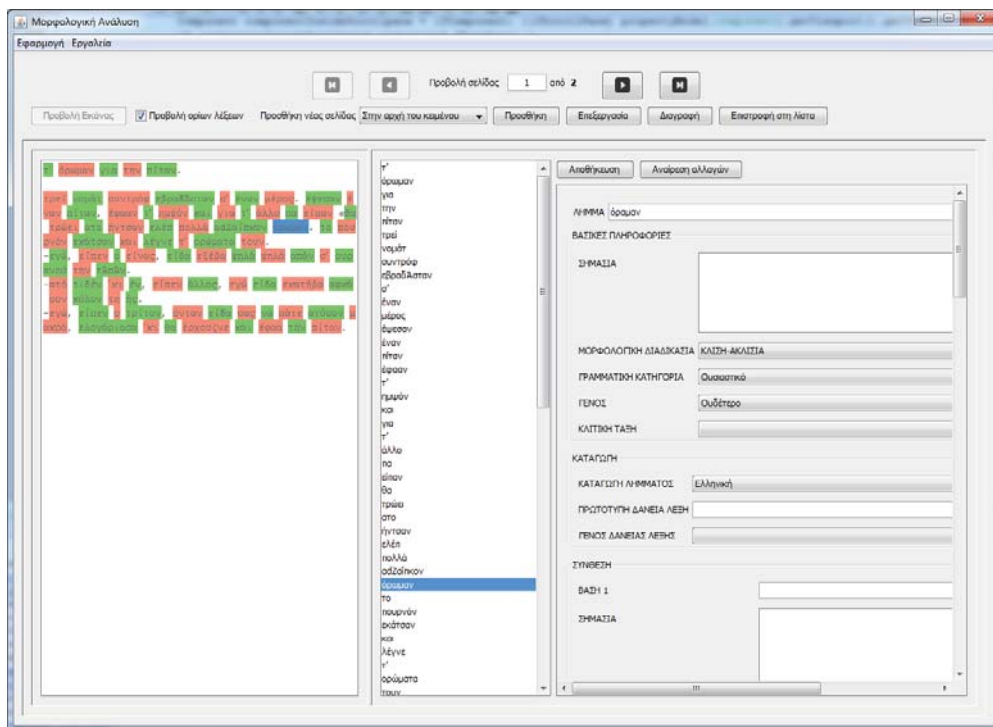
Για την επισημείωση των δύο σωμάτων ακολουθήθηκαν ίδιες στρατηγικές επισημείωσης, τουλάχιστον στα βασικά γλωσσικά επίπεδα. Η ουσιαστικότερη διαφοροποίηση, εντούτοις, εντοπίζεται στο φωνητικό-φωνολογικό επίπεδο, όπου είναι αναμενόμενα να υπάρχουν διαφορετικά επίπεδα επισημείωσης που θα απουσιάζουν (αναλυτικότερα βλ. Κολιοπούλου, Μαρκόπουλος & Παντελίδης 2015).

#### 4.1.1 Μορφολογικό επίπεδο

Και στα δύο σώματα οι κατηγορίες και υποκατηγορίες μορφολογικής ανάλυσης είναι ίδιες, όπου κυριαρχούν οι λίστες με τις προεπιλεγμένες τιμές στις περισσότερες περιπτώσεις. Οι κατηγορίες ανάλυσης περιέχουν πληροφορίες, όπως λήμμα, μορφολογική διαδικασία, γένος, κλιτική τάξη,

γραμματική κατηγορία, καταγωγή, τύποι βάσεων / μορφημάτων / παραγωγικών προσφυμάτων / κλιτικών προσφυμάτων (ανά γραμματική κατηγορία).

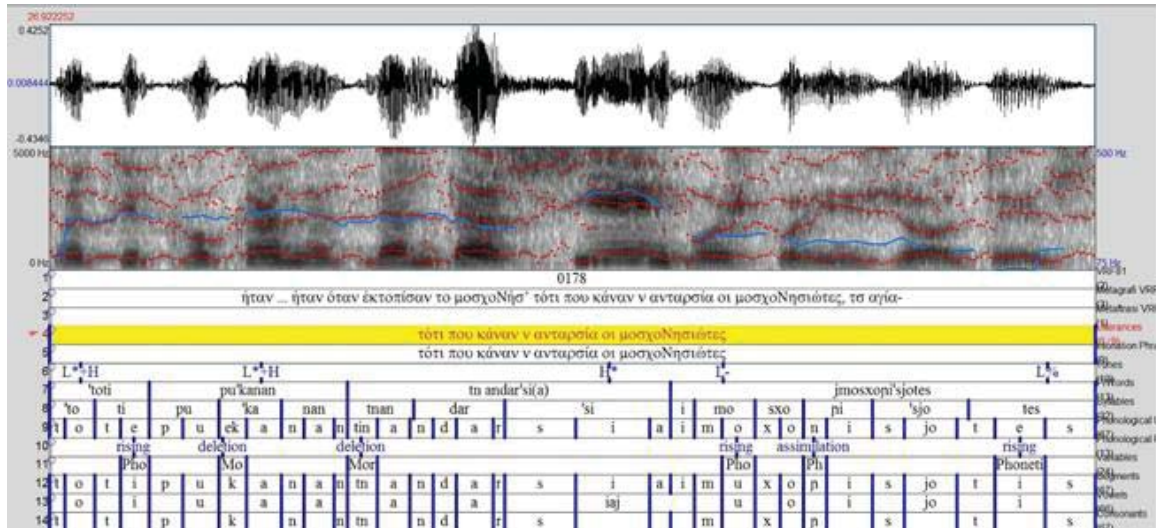
1. ΛΗΜΜΑ	2. ΜΟΡΦΟΛΟΓΙΚΗ ΔΙΑΔΙΚΑΣΙΑ	3. ΓΡΑΜΜΑΤΙΚΗ ΚΑΤΗΓΟΡΙΑ	4. ΓΕΝΟΣ	5. ΚΛΙΤΙΚΗ ΤΑΞΗ	6. ΚΑΤΑΓΩΓΗ ΛΗΜΜΑΤΟΣ
	Κλίση-Ακλίση	Επίθετο	Ουδέτερο	ΚΤ1-Ουσιαστικά	Τουρκική
	Παραγωγή-Κλίση	Ουσιαστικό	Αρσενικό	ΚΤ2-Ουσιαστικά	Ελληνική
	Σύνθεση-Κλίση	Άρθρο	Αρσενικό	ΚΤ3-Ουσιαστικά	Ρομανική
	Σύνθεση-Παραγωγή-Κλίση	Ρήμα	Θηλυκό	ΚΤ4-Ουσιαστικά	Άλλη
	Παραγωγή-Σύνθεση-Κλίση	Επίρρημα	Χωρίς γένος	ΚΤ5-Ουσιαστικά	
		Αντωνυμία		ΚΤ6-Ουσιαστικά	
		Γερούνδιο		ΚΤ7-Ουσιαστικά	
		Απαρέμφατος		ΚΤ8-Ουσιαστικά	
		Μετοχή		ΚΤ9-Ουσιαστικά	
		Επιφώνημα		ΚΤ10-Ουσιαστικά	
		Πρόθεση		ΚΤ1-Ρήματα	
		Σύνδεσμος		ΚΤ2Α-Ρήματα	
		Έκφραση		ΚΤ2Β-Ρήματα	
		Αριθμητικό		ΜΚΤ -Επίθετα	
				Άκλιτο	



Εικόνα 6 & 7: Δείγμα μορφολογικής ανάλυσης στο στάδιο προ-επεξεργασίας και στο στάδιο χρήσης του Morph Tagger

#### 4.1.2 Φωνολογικό-φωνητικό επίπεδο

Η διαφοροποίηση μεταξύ των σωμάτων στο συγκεκριμένο επίπεδο είναι αναμενόμενη. Ενώ στο σώμα γραπτών κειμένων γίνεται εντοπισμός φαινομένων φωνηέντων και συμφώνων (ανάπτυξη, ανομοίωση, αποβολή, ανύψωση, αφομοίωση κτλ.) με μονοεπίπεδο tier, στο σώμα προφορικών κειμένων η πολυεπίπεδη χρήση tiers ανάλυσης συμπεριλαμβάνει ανάλυση εκφωνημάτων, φωνολογικών λέξεων, συλλαβών, φωνημάτων, επιτονισμού, συνεισφορών, κτλ. Γίνεται χρήση μιας τροποποιημένης έκδοσης του IPA για τη συνολική επισήμειωση των ηχητικών αρχείων.



Εικόνα 8: Πολυεπίπεδη φωνολογική επισημείωση διαλεκτικού υλικού στο Praat.

## 4.2 Προκλήσεις στην επισημείωση μεταξύ σωμάτων κειμένων

Η σημαντικότερη πρόκληση και τα σημαντικότερα ερευνητικά ζητήματα εντοπίζονται στην επισημείωση στα σώματα γραπτών κειμένων. Όπως επισημαίνουν οι Κολιοπούλου, Μαρκόπουλος και Παντελίδης (2015) δεν έγινε φωνητική/ φωνολογική μεταγραφή, γιατί:

- (α) Τα ακριβή φωνολογικά χαρακτηριστικά των τριών διαλέκτων παραμένουν αμφίβολα, καθότι τα περισσότερα κείμενα είναι παλαιότερα των 75 ετών και η κωδικοποίηση των φαινομένων έγινε με τυχαίο, μη-επιστημονικό, αλλά συστηματικό τρόπο από τους συγγραφείς.
- (β) Τα γραπτά κείμενα δεν ενδείκνυνται για φωνητική μεταγραφή, γιατί αρκετοί συμβολισμοί δεν μπορούν να αντιστοιχισθούν με σιγουριά στα αντίστοιχα σύμβολα του IPA.
- (γ) Η μη-επιστημικότητα των συντακτών και η αυθαιρεσία συμβόλων στη συγγραφή των κειμένων, εμφανίζεται έντονα στο δείγμα: παράλληλα, βαρύνουσας σημασίας είναι η απουσία συνοδευτικού ενδείκτη ή εισαγωγικού κειμένου που να εξηγούν τις όποιες αποφάσεις πήραν κατά τη συλλογή του υλικού ή τη γραπτή απόδοση των προφορικών μαρτυριών.

Επομένως, για να υπερκεραστούν τα προβλήματα: (α) χρησιμοποιήθηκαν συμπεράσματα από την επισημείωση των προφορικών κειμένων για αμφίβολουσ χαρακτήρες, μιας και τα ηχητικά αρχεία ενδείκνυνται για τέτοια μεταγραφή, (β) έγινε επιβεβαίωση συμβόλων από άλλα κείμενα ίδιας περιόδου, όσο αυτό ήταν εφικτό και (γ) ακολουθήθηκε η χρήση ελληνικού αλφαβήτου με την καθιερωμένη ιστορική ορθογραφία. Στο τελικό έλεγχο των επισημειώσεων τα δύο σώματα κειμένων θα λειτουργήσουν ως ελεγκτές ακρίβειας και συνέπειας για την επικαιροποίηση των προβληματικών επισημειώσεων. Ταυτόχρονα αποτελούν ένα αξιόπιστο δείγμα για επαλήθευση των πινάκων αντιστοίχισης συμβόλων με το IPA.

## 5. Συμπεράσματα

Η συνέπεια στην επισημείωση σωμάτων κειμένων παραμένει σοβαρό ζήτημα για τη διακειμενική γλωσσολογία. Σημαντικά ζητήματα για την επισημείωση σε φωνολογικό επίπεδο αντιμετωπίστηκαν κατά τη μελέτη καθότι έγινε μια συστηματική προσπάθεια να ενοποιηθούν όλες οι διαφορετικές μεταγραφές διαλεκτικού γραπτού υλικού μιας και δεν υπήρχε προηγουμένως κοινή στρατηγική απεικόνισης. Παράλληλα προτείνεται πολυεπίπεδη φωνολογική (παράλληλα με τη μορφολογική) επισημείωση του σώματος κειμένων καθιερώνοντας ένα βασικό πρότυπο επισημείωσης διαλεκτικού

υλικού για τις Νεοελληνικές Διαλέκτους σε καθιερωμένα λογισμικά ανάλυσης ομιλίας, ενώ γίνεται η χρήση των επισημειώσεων για δια-σωματική επικαιροποίηση της συνέπειας και της ακρίβειας της συνολικής επισημείωσης.

## Βιβλιογραφία

- Dickinson, M & W. Detmar-Meurers (2003) Detecting errors in part-of-speech annotation. Στο *Proceedings of EACL-03*. Budapest: Association for Computational Linguistics, 107-114.
- Dickinson, M. (2005) *Error detection and correction in annotated corpora*. Ph.D. thesis. Columbus, Ohio: The Ohio State University.
- Eskin, E. (2000) Automatic corpus correction with anomaly detection. Στο: *Proceedings of NAACL-00*. Seattle: WA, 148-153.
- Gries, S. Th. & A. L. Berez (υπό έκδοση) Linguistic annotation in/for corpus linguistics. Στο: N. Ide & J. Pustejovsky (επιμ.), *Handbook of linguistic annotation*. Berlin & New York: Springer.
- Karanikolas, N. Galiotou, E. & A. Ralli (2014) Towards a unified exploitation of electronic dialectal corpora: problems and perspectives. Στο: P. Sojka et al. (επιμ.), *TSD 2014*. Switzerland: Springer, 257-266.
- Karasimos, A., Melissaropoulou, D., Ralli, A., Papazachariou, D. & D. Asimakopoulos (2008) GREED: Cataloguing and encoding Modern Greek dialectal spoken corpora. Paper presented in *CatCod 2008*, 4-5 December, Orleans, France.
- Ule, T. & K. Simon (2004) Unexpected productions may well be errors. Στο: *Proceeding of LREC 2004*. Lisbon, 1795-1798.
- Van Halteren, H. (2000) The detection of inconsistency in manually tagged text. Στο: A. Abeillé, T. Brants & H. Uszkoreit (επιμ.), *Proceedings of LINC-00*. Luxembourg: MIT-LINC, 48-55.
- Wallis, S.A. & G. Nelson (2001) Knowledge discovery in grammatically analysed corpora. *Data mining and knowledge discovery 15*: 307-340.
- Wells, J. C. (1997) SAMPA computer readable phonetic alphabet. Στο: D. Gibbon, R. Moore & R. Winski (επιμ.), *Handbook of standards and resources for spoken language systems*. Berlin & New York: Mouton de Gruyter, Part IV, Section B.
- Κολιοπούλου, Μ., Μαρκόπουλος, Θ. & Ν. Παντελίδης (2015) Πόντος, Καπαδοκία, Αϊβαλί: προκλήσεις ενός ψηφιακού σώματος γραπτού υλικού. Στο: G. Kotzoglou et al. (επιμ.), *Proceedings of ICGL11*. Rhodes: University of Aegean, 750-759.
- Ράλλη Α., Παπαζαχαρίου Δ. & Α. Καρασίμος (2010) Εργαστήριο Νεοελληνικών Διαλέκτων και η βάση δεδομένων GREED. Στο: A. Ralli et al. (επιμ.), *Proceedings of 4<sup>th</sup> International Conference of Modern Greek Dialects and Linguistic Theory*. Patras: University of Patras, 7-14.
- Μανωλέσσου, Ι., Μπέης, Σ. & Χ. Μπασέα (2012) Η φωνητική μεταγραφή των Νεοελληνικών Διαλέκτων. *Λεξικογραφικόν δελτίον 26*: 161-222.