



FACHLITERATUR
EDITION ROMIOSINI
ΕΝΙΣΤΗΜΗ



ICGL12

12th INTERNATIONAL CONFERENCE
ON GREEK LINGUISTICS
16 – 19 SEPTEMBER 2015
FREIE UNIVERSITÄT BERLIN, CEMOG

Proceedings of the ICGL12

vol. 1

The International Conference on Greek Linguistics is a biennial meeting on the study and analysis of Greek (Ancient, Medieval and Modern), placing particular emphasis on the later stages of the language.

PROCEEDINGS OF THE ICGL12
ΠΡΑΚΤΙΚΑ ΤΟΥ ICGL12

**Thanasis Georgakopoulos, Theodossia-Soula Pavlidou, Miltos Pechlivanos,
Artemis Alexiadou, Jannis Androutsopoulos, Alexis Kalokairinos,
Stavros Skopeteas, Katerina Stathi (Eds.)**

**PROCEEDINGS OF THE 12TH INTERNATIONAL
CONFERENCE ON GREEK LINGUISTICS**

**ΠΡΑΚΤΙΚΑ ΤΟΥ 12^{ΟΥ} ΣΥΝΕΔΡΙΟΥ ΕΛΛΗΝΙΚΗΣ
ΓΛΩΣΣΟΛΟΓΙΑΣ**

VOL. 1

© 2017 Edition Romiosini/CeMoG, Freie Universität Berlin. Alle Rechte vorbehalten.
Vertrieb und Gesamtherstellung: Epubli (www.epubli.de)
Satz und Layout: Rea Papamichail / Center für Digitale Systeme, Freie Universität Berlin
Gesetzt aus Minion Pro
Umschlaggestaltung: Thanasis Georgiou, Yorgos Konstantinou
Umschlagillustration: Yorgos Konstantinou

ISBN 978-3-946142-34-8
Printed in Germany

Online-Bibliothek der Edition Romiosini:
www.edition-romiosini.de

Στη μνήμη του *Gaberell Drachman* (†10.9.2014)
και της Αγγελικής Μαλικούτη-*Drachman* (†4.5.2015)
για την τεράστια προσφορά τους στην ελληνική γλωσσολογία
και την αγάπη τους για την ελληνική γλώσσα

ΣΗΜΕΙΩΜΑ ΕΚΔΟΤΩΝ

Το 12ο Διεθνές Συνέδριο Ελληνικής Γλωσσολογίας (International Conference on Greek Linguistics/ICGL12) πραγματοποιήθηκε στο Κέντρο Νέου Ελληνισμού του Ελεύθερου Πανεπιστημίου του Βερολίνου (Centrum Modernes Griechenland, Freie Universität Berlin) στις 16-19 Σεπτεμβρίου 2015 με τη συμμετοχή περίπου τετρακοσίων συνέδρων απ' όλον τον κόσμο.

Την Επιστημονική Επιτροπή του ICGL12 στελέχωσαν οι Θανάσης Γεωργακόπουλος, Θεοδοσία-Σούλα Παυλίδου, Μίλτος Πεχλιβάνος, Άρτεμις Αλεξιάδου, Δώρα Αλεξοπούλου, Γιάννης Ανδρουτσόπουλος, Αμαλία Αρβανίτη, Σταύρος Ασημακόπουλος, Αλεξάνδρα Γεωργακοπούλου, Κλεάνθης Γκρώμαν, Σαβίνα Ιατρίδου, Mark Janse, Brian Joseph, Αλέξης Καλοκαιρινός, Ναπολέων Κάτσος, Ευαγγελία Κορδώνη, Αμαλία Μόζερ, Ελένη Μπουτουλούση, Κική Νικηφορίδου, Αγγελική Ράλλη, Άννα Ρούσου, Αθηνά Σιούπη, Σταύρος Σκοπετέας, Κατερίνα Στάθη, Μελίτα Σταύρου, Αρχόντω Τερζή, Νίνα Τοπιντζή, Ιάνθη Τσιμπλή και Σταυρούλα Τσιπλάκου.

Την Οργανωτική Επιτροπή του ICGL12 στελέχωσαν οι Θανάσης Γεωργακόπουλος, Αλέξης Καλοκαιρινός, Κώστας Κοσμάς, Θεοδοσία-Σούλα Παυλίδου και Μίλτος Πεχλιβάνος.

Οι δύο τόμοι των πρακτικών του συνεδρίου είναι προϊόν της εργασίας της Εκδοτικής Επιτροπής στην οποία συμμετείχαν οι Θανάσης Γεωργακόπουλος, Θεοδοσία-Σούλα Παυλίδου, Μίλτος Πεχλιβάνος, Άρτεμις Αλεξιάδου, Γιάννης Ανδρουτσόπουλος, Αλέξης Καλοκαιρινός, Σταύρος Σκοπετέας και Κατερίνα Στάθη.

Παρότι στο συνέδριο οι ανακοινώσεις είχαν ταξινομηθεί σύμφωνα με θεματικούς άξονες, τα κείμενα των ανακοινώσεων παρατίθενται σε αλφαβητική σειρά, σύμφωνα με το λατινικό αλφάβητο· εξαίρεση αποτελούν οι εναρκτήριες ομιλίες, οι οποίες βρίσκονται στην αρχή του πρώτου τόμου.

Η Οργανωτική Επιτροπή του ICGL12

ΠΕΡΙΕΧΟΜΕΝΑ

Σημείωμα εκδοτών	7
Περιεχόμενα	9
 Peter Mackridge: <i>Some literary representations of spoken Greek before nationalism(1750-1801)</i>	17
Μαρία Σηφianού: <i>Η έννοια της ευγένειας στα Ελληνικά</i>	45
Σπυριδούλα Βαρλοκώστα: <i>Syntactic comprehension in aphasia and its relationship to working memory deficits</i>	75
 Ευαγγελία Αχλάδη, Αγγελική Δούρη, Ευγενία Μαλικούτη & Χρυσάνθη Παρασχάκη-Μπαράν: <i>Γλωσσικά λάθη τουρκόφωνων μαθητών της Ελληνικής ως ξένης/δεύτερης γλώσσας: Ανάλυση και διδακτική αξιοποίηση</i>	109
Κατερίνα Αλεξανδρή: <i>Η μορφή και η σημασία της διαβάθμισης στα επίθετα που δηλώνουν χρώμα</i>	125
Eva Anastasi, Ageliki Logotheti, Stavri Panayiotou, Marilena Serafim & Charalambos Themistocleous: <i>A Study of Standard Modern Greek and Cypriot Greek Stop Consonants: Preliminary Findings</i>	141
Anna Anastassiadis-Symeonidis, Elisavet Kiourti & Maria Mitsiaki: <i>Inflectional Morphology at the service of Lexicography: ΚΟΜΟΛεξ, A Cypriot Morphological Dictionary</i>	157

Γεωργία Ανδρέου & Ματίνα Τασιούδη: <i>Η ανάπτυξη του λεξιλογίου σε παιδιά με Σύνδρομο Απνοιών στον Ύπνο</i>	175
Ανθούλα- Ελευθερία Ανδρεσάκη: <i>Ιατρικές μεταφορές στον δημοσιογραφικό λόγο της κρίσης: Η οπτική γωνία των Γερμανών</i>	187
Μαρία Ανδριά: <i>Προσεγγίζοντας θέματα Διαγλωσσικής Επίδρασης μέσα από το πλαίσιο της Γνωσιακής Γλωσσολογίας: ένα παράδειγμα από την κατάκτηση της Ελληνικής ως Γ2</i>	199
Spyros Armostis & Kakia Petinou: <i>Mastering word-initial syllable onsets by Cypriot Greek toddlers with and without early language delay</i>	215
Julia Bacsakai-Atkari: <i>Ambiguity and the Internal Structure of Comparative Complements in Greek</i>	231
Costas Canakis: <i>Talking about same-sex parenthood in contemporary Greece: Dynamic categorization and indexicality</i>	243
Michael Chiou: <i>The pragmatics of future tense in Greek</i>	257
Maria Chondrogianni: <i>The Pragmatics of the Modern Greek Segmental Markers</i>	269
Katerina Christopoulou, George J. Xydopoulos & Anastasios Tsangalidis: <i>Grammatical gender and offensiveness in Modern Greek slang vocabulary</i>	291
Aggeliki Fotopoulou, Vasiliki Foufi, Tita Kyriacopoulou & Claude Martineau: <i>Extraction of complex text segments in Modern Greek</i>	307
Αγγελική Φωτοπούλου & Βούλα Γιούλη: <i>Από την «Έκφραση» στο «Πολύτροπο»: σχεδιασμός και οργάνωση ενός εννοιολογικού λεξικού</i>	327
Marianthi Georgalidou, Sofia Lampropoulou, Maria Gasouka, Apostolos Kostas & Xanthippi Foulidi: <i>“Learn grammar”: Sexist language and ideology in a corpus of Greek Public Documents</i>	341
Maria Giagkou, Giorgos Fragakakis, Dimitris Pappas & Harris Papageorgiou: <i>Feature extraction and analysis in Greek L2 texts in view of automatic labeling for proficiency levels</i>	357

Dionysis Goutsos, Georgia Fragaki, Irene Florou, Vasiliki Kakousi & Paraskevi Savvidou: <i>The Diachronic Corpus of Greek of the 20th century: Design and compilation</i>	369
Kleanthes K. Grohmann & Maria Kambanaros: <i>Bilectalism, Comparative Bilingualism, and the Gradience of Multilingualism: A View from Cyprus</i>	383
Günther S. Henrich: „Γεωγραφία νεωτερική“ στο Λίβιστρος και Ροδάμνη: μετατόπιση ονομάτων βαλτικών χωρών προς την Ανατολή;.....	397
Noriyo Hoozawa-Arkenau & Christos Karvounis: <i>Vergleichende Diglossie - Aspekte im Japanischen und Neugriechischen: Veritäten - Interferenz</i>	405
Μαρία Ιακώβου, Ηριάννα Βασιλειάδη-Λιναρδάκη, Φλώρα Βλάχου, Όλγα Δήμα, Μαρία Καββαδία, Τατιάνα Κατσίνα, Μαρίνα Κουτσουμπού, Σοφία-Νεφέλη Κύτρου, Χριστίνα Κωστάκου, Φρόσω Παππά & Σταυριαλένα Περρέα: <i>ΣΕΠΙΜΕ2: Μια καινούρια πηγή αναφοράς για την Ελληνική ως Γ2</i>	419
Μαρία Ιακώβου & Θωμάς Ρουσουλιώτη: <i>Βασικές αρχές σχεδιασμού και ανάπτυξης του νέου μοντέλου αναλυτικών προγραμμάτων για τη διδασκαλία της Ελληνικής ως δεύτερης/ξένης γλώσσας</i>	433
Μαρία Καμηλάκη: «Μαζί μου ασχολείσαι, πόσο μαλάκας είσαι!»: Λέξεις-ταμπού και κοινωνιογλωσσικές ταυτότητες στο σύγχρονο ελληνόφωνο τραγούδι.....	449
Μαρία Καμηλάκη, Γεωργία Κατσούδα & Μαρία Βραχιονίδου: <i>Η εννοιολογική μεταφορά σε λέξεις-ταμπού της ΝΕΚ και των νεοελληνικών διαλέκτων</i>	465
Eleni Karantzola, Georgios Mikros & Anastassios Papaioannou: <i>Lexico-grammatical variation and stylometric profile of autograph texts in Early Modern Greek</i>	479
Sviatlana Karpava, Maria Kambanaros & Kleanthes K. Grohmann: <i>Narrative Abilities: MAINing Russian–Greek Bilingual Children in Cyprus</i>	493
Χρήστος Καρβούνης: <i>Γλωσσικός εξαρχαϊσμός και «ιδεολογική» νόρμα: Ζητήματα γλωσσικής διαχείρισης στη νέα ελληνική</i>	507

Demetra Katis & Kiki Nikiforidou: <i>Spatial prepositions in early child Greek: Implications for acquisition, polysemy and historical change</i>	525
Γεωργία Κατσούδα: <i>Το επίθημα -ούνα στη ΝΕΚ και στις νεοελληνικές διαλέκτους και ιδιώματα</i>	539
George Kotzoglou: <i>Sub-extraction from subjects in Greek: Its existence, its locus and an open issue</i>	555
Veranna Kyprioti: <i>Narrative, identity and age: the case of the bilingual in Greek and Turkish Muslim community of Rhodes, Greece</i>	571
Χριστίνα Λύκου: <i>Η Ελλάδα στην Ευρώπη της κρίσης: Αναπαραστάσεις στον ελληνικό δημοσιογραφικό λόγο</i>	583
Nikos Liosis: <i>Systems in disruption: Propontis Tsakonian</i>	599
Katerina Magdou, Sam Featherston: <i>Resumptive Pronouns can be more acceptable than gaps: Experimental evidence from Greek</i>	613
Maria Margarita Makri: <i>Opos identity comparatives in Greek: an experimental investigation</i>	629
2ος Τόμος	
Περιεχόμενα	651
Vasiliki Makri: <i>Gender assignment to Romance loans in Katoitaliótika: a case study of contact morphology</i>	659
Evgenia Malikouti: <i>Usage Labels of Turkish Loanwords in three Modern Greek Dictionaries</i>	675
Persephone Mamoukari & Penelope Kambakis-Vougiouklis: <i>Frequency and Effectiveness of Strategy Use in SILL questionnaire using an Innovative Electronic Application</i>	693

Georgia Maniati, Voula Gotsoulia & Stella Markantonatou: <i>Contrasting the Conceptual Lexicon of ILSP (CL-ILSP) with major lexicographic examples</i>	709
Γεώργιος Μαρκόπουλος & Αθανάσιος Καρασίμος: <i>Πολυεπίπεδη επισημείωση του Ελληνικού Σώματος Κειμένων Αφασικού Λόγου</i>	725
Πωλίνα Μεσηνιώτη, Κατερίνα Πούλιου & Χριστόφορος Σουγανίδης: <i>Μορφοσυντακτικά λάθη μαθητών Τάξεων Υποδοχής που διδάσκονται την Ελληνική ως Γ2</i>	741
Stamatia Michalopoulou: <i>Third Language Acquisition. The Pro-Drop-Parameter in the Interlanguage of Greek students of German</i>	759
Vicky Nanousi & Arhonto Terzi: <i>Non-canonical sentences in agrammatism: the case of Greek passives</i>	773
Καλομοίρα Νικολού, Μαρία Ξεφτέρη & Νίτσα Παραχεράκη: <i>Το φαινόμενο της σύνθεσης λέξεων στην κυκλαδοκρητική διαλεκτική ομάδα</i>	789
Ελένη Παπαδάμου & Δώρης Κ. Κυριαζής: <i>Μορφές διαβαθμιστικής αναδίπλωσης στην ελληνική και στις άλλες βαλκανικές γλώσσες</i>	807
Γεράσιμος Σοφοκλής Παπαδόπουλος: <i>Το δίπολο «Εμείς και οι Άλλοι» σε σχόλια αναγνωστών της Lifo σχετικά με τη Χρυσή Αυγή</i>	823
Ελένη Παπαδοπούλου: <i>Η συνδυαστικότητα υποκοριστικών επιθημάτων με β' συνθετικό το επίθημα -άκι στον διαλεκτικό λόγο</i>	839
Στέλιος Πιπερίδης, Πένυ Λαμπροπούλου & Μαρία Γαβριηλίδου: <i>clarin:el. Υποδομή τεκμηρίωσης, διαμοιρασμού και επεξεργασίας γλωσσικών δεδομένων</i>	851
Maria Pontiki: <i>Opinion Mining and Target Extraction in Greek Review Texts</i>	871
Anna Roussou: <i>The duality of mipos</i>	885

Stathis Selimis & Demetra Katis: <i>Reference to static space in Greek: A cross-linguistic and developmental perspective of poster descriptions</i>	897
Evi Sifaki & George Tsoulas: <i>XP-V orders in Greek</i>	911
Konstantinos Sipitanos: <i>On desiderative constructions in Naousa dialect</i>	923
Eleni Staraki: <i>Future in Greek: A Degree Expression</i>	935
Χριστίνα Τακούδα & Ευανθία Παπαευθυμίου: <i>Συγκριτικές διδακτικές πρακτικές στη διδασκαλία της ελληνικής ως Γ2: από την κριτική παρατήρηση στην αναπλαισίωση</i>	945
Alexandros Tantos, Giorgos Chatziioannidis, Katerina Lykou, Meropi Papatheohari, Antonia Samara & Kostas Vlachos: <i>Corpus C58 and the interface between intra- and inter-sentential linguistic information</i>	961
Arhonto Terzi & Vina Tsakali: <i>The contribution of Greek SE in the development of locatives</i>	977
Paraskevi Thomou: <i>Conceptual and lexical aspects influencing metaphor realization in Modern Greek</i>	993
Nina Topintzi & Stuart Davis: <i>Features and Asymmetries of Edge Geminates</i>	1007
Liana Tronci: <i>At the lexicon-syntax interface Ancient Greek constructions with ἔχειν and psychological nouns</i>	1021
Βίλλυ Τσάκωνα: <i>«Δημοκρατία είναι 4 λύκοι και 1 πρόβατο να ψηφίζουν για φαγητό»:Αναλύοντας τα ανέκδοτα για τους/τις πολιτικούς στην οικονομική κρίση</i>	1035
Ειρήνη Τσαμαδού- Jacobberger & Μαρία Ζέρβα: <i>Εκμάθηση ελληνικών στο Πανεπιστήμιο Στρασβούργου: κίνητρα και αναπαραστάσεις</i> ...	1051
Stavroula Tsiplakou & Spyros Armotistis: <i>Do dialect variants (mis)behave? Evidence from the Cypriot Greek koine</i>	1065
Αγγελική Τσόκογλου & Σύλα Κλειδή: <i>Συζητώντας τις δομές σε -οντας</i>	1077

Αλεξιάννα Τσότσου:

Η μεθοδολογική προσέγγιση της εικόνας της Γερμανίας στις ελληνικές εφημερίδες 1095

Anastasia Tzilinis:

Begründendes Handeln im neugriechischen Wissenschaftlichen Artikel: Die Situierung des eigenen Beitrags im Forschungszusammenhang..... 1109

Κυριακούλα Τζωρτζάτου, Αργύρης Αρχάκης, Άννα Ιορδανίδου & Γιώργος Ι. Ξυδόπουλος:
Στάσεις απέναντι στην ορθογραφία της Κοινής Νέας Ελληνικής: Ζητήματα ερευνητικού σχεδιασμού 1123

Nicole Vassalou, Dimitris Papazachariou & Mark Janse:

The Vowel System of Mišótika Cappadocian 1139

Marina Vassiliou, Angelos Georgaras, Prokopis Prokopidis & Haris Papageorgiou:

Co-referring or not co-referring? Answer the question!..... 1155

Jeroen Vis:

The acquisition of Ancient Greek vocabulary..... 1171

Christos Vlachos:

Mod(aliti)es of lifting wh-questions..... 1187

Ευαγγελία Βλάχου & Κατερίνα Φραντζή:

Μελέτη της χρήσης των ποσοδεικτών λίγο-λιγάκι σε κείμενα πολιτικού λόγου 1201

Madeleine Voga:

Τι μας διδάσκουν τα ρήματα της ΝΕ σχετικά με την επεξεργασία της μορφολογίας 1213

Werner Voigt:

«Σεληνάκι μου λαμπρό, φέγγε μου να περπατώ ...» oder: warum es in dem bekannten Lied nicht so, sondern eben φεγγαράκι heißt und ngr. φεγγάρι

1227

Μαρία Βραχιονίδου:

Υποκοριστικά επιρρήματα σε νεοελληνικές διαλέκτους και ιδιώματα 1241

Jeroen van de Weijer & Marina Tzakosta:

*The Status of *Complex in Greek.....* 1259

Theodoros Xioufis:

The pattern of the metaphor within metonymy in the figurative language of romantic love in modern Greek..... 1275

THE DIACHRONIC CORPUS OF GREEK OF THE 20TH CENTURY: DESIGN AND COMPILATION

Dionysis Goutsos, Georgia Fragaki, Irene Florou,
Vasiliki Kakousi & Paraskevi Savvidou

National and Kapodistrian University of Athens

dgoutsos@phil.uoa.gr, efraga@phil.uoa.gr, eirini.florou@gmail.com, vaswk1202@hotmail.com & psavvidou@phil.uoa.gr

Περίληψη

Στο άρθρο παρουσιάζεται το Διαχρονικό Σώμα Ελληνικών Κειμένων του 20ού αιώνα, το πρώτο διαχρονικό σώμα κειμένων της ελληνικής, που έχει σχεδιαστεί για τη μελέτη της πρόσφατης γλωσσικής αλλαγής στα ελληνικά. Ειδικότερα, παρουσιάζονται ζητήματα που αφορούν τη συλλογή γλωσσικών δεδομένων του εικοστού αιώνα στα ελληνικά, η σύνθεση του σώματος κειμένων (γένη, είδη, αριθμός λέξεων κ.λπ.), τα παραδοτέα του ερευνητικού προγράμματος που οδήγησε στη δημιουργία του σώματος κειμένων, καθώς και ορισμένα προκαταρκτικά ευρήματα από την ανάλυσή του.

Keywords: corpus design and compilation, diachronic corpus, recent language change

1. Diachronic corpora and Greek

Corpus linguistics has considerably improved the description of languages by allowing access to large bodies of authentic texts, as well as by contributing to a broad range of applications in lexicography, the writing of grammars, lexical semantics, language teaching, the study of language and ideology, translation, media studies etc. (see, among else, Hunston 2002: 13-14, Meyer 2002: 1-29, Baker et al. 2006). Unlike other

languages, Greek has only benefited to a small extent by the development of this field, mainly because of the lack of large Greek corpora, with the exception of the *Hellenic National Corpus* (HNC, 47 million words, texts published from 1976 to 2007) and the *Corpus of Greek Texts* (CGT, 30 million words, texts from 1990 to 2010).¹ Both can be characterized as synchronic corpora, in the sense that they offer a view of a specific period of the Greek language.

This paper presents the design and compilation of the *Diachronic Corpus of Greek of the 20th century* ((Greek Corpus 20 or GC20)), the first diachronic corpus of Greek, developed with a view of studying recent language change.² Its goal is to gather 20 million words from Greek texts coming from the first nine decades of the 20th century, to be integrated with the existing 30 million word CGT, which includes texts from the 1990s onwards.

Historical or diachronic corpora have been compiled or are under preparation for other languages or language varieties like the *Helsinki Corpus of English Texts*, which covers Old, Middle and Early Modern English, the *Corpus of Historical English Registers* (ARCHER), which contains British and American English texts from 1650 to the present, the four corpora including *Brown* and *Frown*, *LOB* and *FLOB*, which can together supply evidence for change in the two varieties of English between 1961 and 1991-1992, *DiaCoris* for Italian etc.³ In practice, three types of corpora have been used to study recent language change in most languages:

- a) diachronic corpora, e.g. the *Corpus of Historical American English* (COHA, with data from 1810 to 2009),
- b) corpus families, e.g. the *Lancaster-Oslo-Bergen Corpus*, with data from every 30 years in the 20th century, including *BLOB-1931* (1928-1934), *LOB* (1961) and *F-LOB* (1991),
- c) synchronic, monitor corpora, e.g. the *British National Corpus* (BNC), including data from 1960 up to now and thus offering a large time span of linguistic material.

1 For more details, see Hatzigeorgiu et al. (2001) for HNC and Goutsos (2010) for CGT.

2 For a definition of recent language change, see Mair (2009: 1120), Davies (2011, 2012).

3 For more details on existing diachronic corpora, see Onelli et al. (2006), Beal et al. (2007), Mair (2009), Baker (2010: 57 ff.), Partington (2010), Aarts et al. (2013).

Diachronic corpora are also of different sizes, from the very big (e.g. *Corpus Diacrónico del Español* with 125 million words), to big (e.g. *O corpus do Português*, with 45 million words), medium-sized (e.g. *Diachronic Czech National Corpus* with 2 million words, *Helsinki Corpus of English Texts* with 1,5 million words) and small (e.g. *Sheffield Corpus of Chinese*, with 18.000 words).

It is also important to notice that a wide range of linguistic phenomena have been studied in diachronic corpora, including vocabulary changes (Baker 2011), grammatical change (Leech et al. 2009), diachronic morphological processes (Baayen & Renouf 1996, Fischer 1998, Duguid 2010), development of phraseology (Davies 2012) and cultural changes (Baker 2010, Marchi 2010, Partington 2012).

Greek has not had a similar diachronic corpus for a number of reasons, among which, as we will discuss below, the difficulty of collecting data is surely prominent. Extra-linguistic factors, such as the socio-historical background in Greece of the 20th century, can account for the lack of data or the occurrence of minimal data for several periods. In addition, linguistic factors such as the persisting diglossia, which is related with important socio-historical events throughout the 20th century, complicate issues of data collection and analysis. For this reason, the development of a diachronic corpus for Greek of the 20th century has been more than imperative.

The research project for the development of GC20 has had the following aims:⁴

- a) to examine the issues involved in the compilation of a diachronic corpus of Greek of the 20th century, including the availability of data across decades, the availability and continuity of text types, and the issue of representativeness;
- b) on the basis of exploration of data sources, to collect data for a diachronic corpus of Greek of the 20th century;
- c) to analyze the corpus with a view to drawing basic conclusions on linguistic change across the decades of the 20th century.

4 The project was funded in the frame of the action “Aristeia I” (“Excellence I”) by the European Cohesion Fund and the Greek government (General Secretariat for Research and Technology). We are grateful to the project’s international advisory board, namely Claudia Claridge (University of Duisburg-Essen), Mark Davies (Brigham Young University), Hendrik De Smet (KU Leuven), Susan M. Fitzmaurice (University of Sheffield), Marianne Hundt (University of Zurich), Christian Mair (University of Freiburg), Terttu Nevalainen (University of Helsinki), Fabio Tamburini (Università di Bologna) and Sean Wallis (University College London), for their help in the various stages of the project.

With respect to these aims, in what follows we first outline the most important issues concerning the compilation of the corpus, we then present the data collected in the corpus and, finally, we discuss the project's deliverables and some preliminary findings.

2. Issues regarding the compilation of the corpus

The project, first, investigated the availability of data in different text types, the feasibility of collecting particular data categories and the possibility of collecting as much data as possible. A major problem concerning the collection of Greek data of the 20th century concerns, first of all, the lack of fully functioning OCR processing facilities for polytonic Greek, the script in which Greek was written for most of the 20th century (specifically, up to 1982). We have developed our own tools by training the open source OCR engine *Tesseract*⁵ with Greek polytonic data and have created a platform, which will be freely available to researchers after the end of the project. However, extensive training is still needed for a fully satisfactory processing of polytonic texts and thus post-editing for several genres has been time-consuming with the effect that it was not possible to process more data. It is expected that, once this platform is available, further training on Greek polytonic data will be easier.

Furthermore, the lack of freely available archives for many Greek genres has been a serious obstacle in data collection. Specifically, Greek TV archives, after a two year period of sudden closure (2013-2015) have become publicly available again, but do not keep news data. In addition, public radio archives are not publicly available. Parliament proceedings are only available online at the site of the Hellenic Parliament for 1900-1935 and from the end of 1989 onwards, leaving thus out five decades of the 20th century. Newspaper archives, especially those of major newspapers that were published for most of the 20th century (e.g. *Kathimerini*, *Vima*) have limited or no access and, despite our efforts to gain access, no progress has been made.

Most importantly, archives that were made open access in the 1990s and 2000s mainly keep image rather than OCR-processed records with the effect that further processing is needed. A notorious example of this practice concerns the online archives of the influential 20th century journal *Nea Estia*, which cannot be processed by any means, but can only be leafed through like a hard copy. Another example concerns

5 The software is available at: <https://github.com/tesseract-ocr/tesseract>

the newspaper archives of the National Library of Greece, which have been processed by a shallow OCR engine, but for which fully OCR-processed files are missing. This problem mostly affects genres of journalistic texts, which constitute a large part of modern Greek synchronic corpora (Goutsos 2010), as well as public records of spoken material, which are sadly underdeveloped for Greek.

A third major issue has to do with the continuity of text types, i.e. the fact that several text types may only be found in certain decades. This is a well-known problem in the diachronic corpora literature (see e.g. Nevalainen & Raumolin-Brunberg 2003: 28) and particularly affects Greek 20th century data in major genres like popularized non-fiction texts. For instance, although there have been several literary journals in the 20th century, no magazines on other subjects seem to be easily accessible for the whole of this century. This is partly an effect of digitization policies, which have exclusively focused on literary journals, especially for the 19th century and the beginning of the 20th century (e.g. the University of Patras collections), but also reflects the fluid limits of general interest magazines for the first half of the 20th century, which mostly included literary contributions (Karaoglou 2005).

At the same time, electronic media-related text types emerged quite late in Greece, with sound films and radio stations appearing in the 1930s and public TV in the 1960s. Full operation of these media was further delayed because of the effects of the Second World War in the 1940s and the military dictatorship of the 1970s.

Taking into account these problems and based on our experience from a pilot version, we decided to follow a double strategy consisting in concentrating on a subset of genres to be fully processed and integrated in the final corpus data, while for the other genres it was decided to collect as much data as possible with a view to processing and including them in the future. Specifically, for full processing it was decided to focus on the genres of Spoken News, Public Speech and Conversation, as regards the spoken mode, and Literature, Academic, Popularized Non-Fiction and Private, as regards the written mode. Data were collected but have not been fully processed and integrated for the genres of Interview, as regards the spoken mode, and News, Opinion Articles, Information Items and Procedural, as regards the written mode.⁶ This decision accords well with the trend noted by Nevalainen & Raumolin-Brunberg (2003: 27) to move

6 Sources for these data include, among else, the National Library of Greece, with which there has been an agreement for data sharing, the Greek Parliament Library collections, mainly for newspapers up to the 1930s, and other private collections.

“from textually balanced multi-purpose corpora towards larger single-genre corpora”, although in our case it is based on the idea of developing micro-corpora for text types found only in certain decades as part of the initial corpus design. It must also be noted that the text types that were fully processed and integrated in the final corpus give emphasis on speech-like (private letters), speech-based (public speeches) and speech-purposed (films, drama, newsreels) text types (cf. Culpeper & Kytö 2010). In this sense, the final corpus is oriented towards data that are more likely to reveal actual speaking patterns of the past.

Mode	Genres	Text types	Codes	Number of words
Spoken	Spoken news	Newsreels	SRF01	78,441
	Public speeches	Parliament Academic Other	STL16 SAL06 SOL16	339,194
	Conversation	Film scripts	SFF19	208,207
Written	Literature	Novels Short stories Poetry Drama	WFB08 WFB09 WFC11 WFB12	1,355,629
	Academic	Humanities Social/Finance Science	WAB13 WAB14 WAB15	1,044,200
	Popularized Non-fiction		WLB13 WLB14	285,252
	Law and administration		WDC34 WDC35	265,924
	Private	Letters	WPO26	188,856
	Miscellanea		WMO99	1,136
Total	3,766,839			

Table 1 | Composition of Greek Corpus 20 (data integrated in May 2016)

3. Corpus composition

Table 1 presents in detail the number of words integrated so far for each genre and text type in all decades covered in GC20, as of May 2016. The total number of words integrated in the corpus so far is 3,766,839, which roughly corresponds to 20% of the target for GC20. It is estimated that the data collected for the genres that have not been integrated in the final corpus amount to more than 15 million words, a figure which covers the remaining percentage of the corpus target, although it is hard to be accurate with non-OCR processed texts. In all, the divergence from the projected target is indicative of the problems related with corpus compilation pointed out in the previous section.

Figure 1 presents the distribution of the data that have been integrated in GC20 across the nine decades of the 20th century. Surprisingly enough, more data have been integrated for the first two decades of the 20th century, while there is a slight increase from the 1950s onward. This may reflect the availability of existing data, as most projects collected data have concentrated on the 19th and the beginning of the 20th century, in particular with respect to literary texts and journals, which have been thought to be of special value. (Copyright restrictions also account for the less easy access to more recent material). It is also clear that more effort is needed to collect data from the 1920s, 1930s and the 1950s.

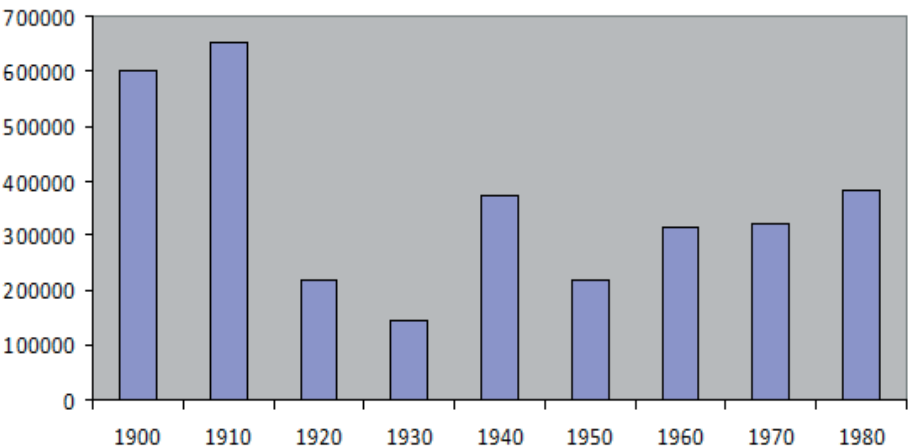


Figure 1 | Current distribution of data in Greek Corpus 20

4. Deliverables and preliminary findings

The project's deliverables include, among else, the compilation of a bibliography on diachronic corpora and an inventory of research projects of diachronic corpora. Both are extensive and offer an updated picture of research conducted in a large variety of languages on issues of language change with the help of diachronic corpora. This is necessary background material for anyone attempting a state-of-the-art description of diachronic corpus research. They also include the oral presentations of an international workshop on the compilation and analysis of diachronic corpora, available online, the various reports, evaluations and publications that were made in the successive stages of the corpus compilation and the project's webpage, which gives access to the corpus itself.⁷

The corpus webpage gives access to all data that has been fully processed, while more texts are constantly being added. The total number of words (tokens) and of their unique occurrence (types) included in the corpus appears in each search. In search one can type in the word or phrase (up to 4 words) they are interested in, using any tone or other diacritic, in order to take all possible versions of the word form occurring in the corpus. For instance, all relevant word forms will appear if you type in *ήμερα* or *ήμερα* or even *ήμερα*. One can also search for part of a word (but not phrase) using wildcards; for instance, the search *ημ*ρα* will give all word forms for *ήμερα* and *ήμετερα* and the search *ημ?ρα* all word forms for *ήμερα*. Results can be sorted according to the node word or phrase, the previous or the next word, the text type in which word forms appear and the data of the texts in which they occur, in ascending or descending order.

Figure 2 presents a screen from a corpus query result, whereas Figure 3 that follows on the second next page illustrates further search statistics provided on the webpage about the frequency development of the words found for the query throughout the nine decades, within a decade etc.

Our preliminary findings from the analysis of the corpus have been reported in Goutsos & Fragaki (2014) and Fragaki & Goutsos (2015) and suggest exciting prospects for the analysis of recent language change in Greek. Thus, an analysis of grammatical words of Greek at the top of the corpus word frequency list, such as *διά* vs. *γιά/για* 'for' or *είς* vs. *σέ/σε* 'in/at' can be revealing of the complex patterns of Greek diglossia in the 20th century. Data analysis supports a variationist view of language

⁷ The corpus is freely available at: <http://greekcorpus20.sek.edu.gr/>



Διαχρονικό σώμα ελληνικών κειμένων του 20ού αιώνα

Αναζήτηση Βοήθεια/Help Ρυθμίσεις Διαχείριση Αποσύνδεση

Αναζήτηση

ΤΑΞΙΝΟΜΗΣΗ	προηγούμενη λέξη ▲ ▼	λέξη/φράση ▲ ▼	επόμενη λέξη ▲ ▼	κειμενικό είδος ▲ ▼	χρονολογία ▲ ▼
1	ἐπιστρέψας εἰς τὸν Β τὰς 3δ. ἡ	ἡμέρα	ἦτο θερμὴ καὶ ἐντὸς ὀλίγου ὁ Β ἐδίψασε	ΑΚΑΔΗΜΑΙΚΑ	1952
2	«Ἢ ἀκριβὴς	ἡμέρα	τοῦ θανάτου μιᾶς μεγάλης δοξασίας εἶνε ἡ <i>ἡμέρα</i> καθ' ἣν ἀρχίζει	ΑΚΑΔΗΜΑΙΚΑ	1963
3	Μουσολίνι ἐδήλωσεν ὅτι «Ἡ 18η Νοεμβρίου (ἡμέρα	ἐγκαθιδρύσεως τοῦ νέου οἰκονομικοῦ	ΑΚΑΔΗΜΑΙΚΑ	1943
4	δὲν ἐβράδυνε νὰ φανῇ. Οὕτω ἦλθε ἡ	ἡμέρα	καθ' ἣν ἡ περὶ τὴν Ἀιττικὴν θάλασσα	ΑΚΑΔΗΜΑΙΚΑ	1949
5	παραδῶση τὴν πόλιν εἰς τὸ πῦρ. ἡ	ἡμέρα	αὐτὴ ὑπῆρξε διὰ τὸν Πειραιᾶ ἡ μελανωτέρα	ΑΚΑΔΗΜΑΙΚΑ	1949
6	σημαντικὴ ἀύξηση βρέθηκε τὴν πρώτη	ἡμέρα	μετὰ τὸν τοκετό, ποῦ μετὰ ἔπεσε σιγὰ	ΑΚΑΔΗΜΑΙΚΑ	1984
7	σὲ φυσιολογικὰ ἐπίπεδα τῆ δεκάτῃ	ἡμέρα	. Βρέθηκε ἐπίσης ὅτι ἡ ποσότης τῆς	ΑΚΑΔΗΜΑΙΚΑ	1984
8	τῆς δραστηριότητος στὸν ὀρό	ἡμέρα	μὲ τὴν <i>ἡμέρα</i> εἶναι κάτι σοβαρότερο ἀπὸ	ΑΚΑΔΗΜΑΙΚΑ	1984
9	τῆς ἀπαιτήσεως τοῦ Τελωνείου. ἡ	ἡμέρα	ἡ ἐνδεικνυομένη ὡς ἀφετηρία τῆς παραγραφῆς	ΑΚΑΔΗΜΑΙΚΑ	1947
10	τῆς ὑποφύσεως. Ἀπὸ τῆς 16ης ὅμως	ἡμέρα	τοῦ κόκλου διαπιστοῦται καταφανῆς	ΑΚΑΔΗΜΑΙΚΑ	1940
11	ἀπ' αὐτὰ τῶν ἐνγλίκων τὴν πέμπτῃ	ἡμέρα	. Ὁ Οκίνακα καὶ ἄλλοι 129 βόηκαν	ΑΚΑΔΗΜΑΙΚΑ	1984
12	ἐπίπεδα ἀπὸ τὴν τρίτῃ ἢ τέταρτῃ	ἡμέρα	. Αὐτὰ τὰ εὐρήματα ἔχουν ἐπανειλημμένα	ΑΚΑΔΗΜΑΙΚΑ	1984
13	ἐπανηγυρίζον τὸ γεγονός ἐνομίζοντες ὅτι ἡ	ἡμέρα	ἐκείνη ἦτο ἡ πρώτη τῆς ἐλληνικῆς	ΑΚΑΔΗΜΑΙΚΑ	1949

Figure 2 | Greek Corpus 20 query result

change on the basis of the thoroughly attested role of frequency (Schneider 2004); specifically, demotic (or Low) variants in Greek diglossia show a U-curve, rather than the expected S-curve of sociolinguistic variation, whereas katharevousa (or High) variants show a “roller-coaster” pattern that is indicative of their stereotypical (in Labov’s sense) or emblematic use. A full-scale investigation of variants like these is expected to contribute to an informed view on standardisation and a better understanding of what happened in the Greek of the 20th century.

Secondly, corpus data support the general principle that recent language change in Greek largely depends on genre (see e.g. Taavitsainen et al. 2015). For example, in film

scripts and literature there is steady preference for Low variants across the century. By contrast, in academic texts and public speeches High variants are preferred in most decades before the 1960s, when there is a sudden rise of Low variants. Newsreels show a haphazard pattern, conforming to the expected rise of Low variants only after the 1960s, whereas private letters are the only genre in which the expected gradual rise of Low variants across all decades is found. This latter finding underlines the importance of collecting and analyzing private letters in understanding recent language change (cf. Dossena & Del Lungo Camiciotti 2012).

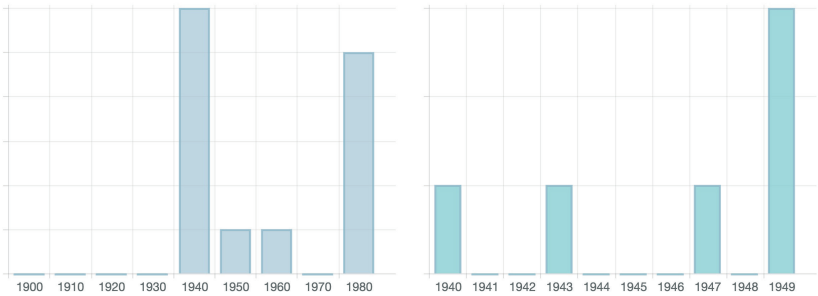
11	της άπαιτήσεως τοῦ Τελωνείου. Ἡ	ἡμέρα	ἡ ἐνδεικνυμένη ὡς ἀφετηρία τῆς παραγραφῆς	ΑΚΑΔΗΜΑΙΚΑ
12	Μουσαλίνι ἐδῆλωσεν ὅτι «Ἡ 18η Νοεμβρίου (ἡμέρα	ἐγκαθιδρύσεως τοῦ νέου οἰκονομικοῦ	ΑΚΑΔΗΜΑΙΚΑ
13	τῆς ἀπορόσεως. Ἀπὸ τῆς 16ης ὁμως	ἡμέρα	τοῦ κύκλου διαπιστοῦται καταφανῆς	ΑΚΑΔΗΜΑΙΚΑ

Στατιστικά αναζήτησης

Αριθμός εμφανίσεων:

ἡμέρα: 15

Αριθμός κειμένων στα οποία απαντά η λέξη/φράση: 7



Σύνολο: 15

Συνολικός αριθμός λέξεων στις οποίες έγινε αναζήτηση: 1982634 (Τύποι: 89978)



Figure 3 | Greek Corpus 20 query statistics

It is clear that the study of sociolinguistic phenomena such as the Greek diglossia will be greatly helped by diachronic corpora such as GC20, which give access to evidence about what actual people said and wrote (language use) rather than what they believed (language attitudes). The analysis of data from the *Diachronic Corpus of Greek of the 20th century* can provide secure indications about the questions surrounding Greek diglossia, by clarifying e.g. whether it is related to the spoken vs. written dichotomy, by identifying when changes took place and by establishing how public attitudes influence the private use of language.

More generally, it is expected that GC20 will offer an invaluable resource for the study of Greek language and culture, providing a point of reference for diachronic research in the still limited spectrum of Greek corpora. Since GC20 was designed to complement the synchronic CGT (Goutsos 2010), the range of available material for Greek is greatly extended. Future perspectives include both the integration of further genres and texts and the morphosyntactic annotation of the corpus, something that has not been attempted before for polytonic Modern Greek.

References

- Aarts, Bas, Close, Joanne, Leech, Geoffrey, and Sean Wallis, eds. 2013. *The Verb Phrase in English: Investigating Recent Language Change with Corpora*. Cambridge: Cambridge University Press.
- Baayen, Harald R. and Antoinette Renouf 1996. "Chronicling the *Times*: Productive lexical innovations in an English newspaper." *Language* 72 (1):69–96.
- Baker, Paul. 2010. *Sociolinguistics and Corpus Linguistics*. Edinburgh: Edinburgh University Press.
- Baker, Paul. 2011. "Times may change, but we will always have money: Diachronic variation in recent British English." *Journal of English Linguistics* 39 (1):65–88.
- Baker, Paul, Hardie, Andrew, and Tony McEnery. 2006. *A Glossary of Corpus Linguistics*. Edinburgh: Edinburgh University Press.
- Beal, Joan C., Corrigan, Karen P., and Hermann L. Moisl, eds. 2007. *Creating and Digitizing Language Corpora. Volume 2: Diachronic Databases*.

- Basingstoke: Palgrave Macmillan.
- Culpeper, Jonathan, and Merja Kytö. 2010. *Early Modern English Dialogues: Spoken Interaction as Writing*. Cambridge: Cambridge University Press.
- Davies, Mark. 2011. "Synchronic and diachronic uses of corpora". In *Perspectives on Corpus Linguistics*, edited by Vander Viana, Sonia Zyngier, and Geoff Barnbrook, 63–80. Amsterdam/Philadelphia: John Benjamins.
- Davies, Mark. 2012. "Examining recent changes in English: Some methodological issues". In *The Oxford Handbook of the History of English*, edited by Terttu Nevalainen, and Elizabeth Closs Traugott, 263–287. Oxford: Oxford University Press.
- Dossena, Marina, and Gabriella Del Lungo Camiciotti. 2012. *Letter Writing in Late Modern Europe*. Amsterdam/Philadelphia: Benjamins.
- Duguid, Alison. 2010. "Investigating anti and some reflections on Modern Diachronic Corpus-Assisted Discourse Studies (MD-CADS)." *Corpora* 5 (2):191–220.
- Fischer, Roswitha. 1998. *Lexical Change in Present-Day English. A Corpus Study of the Motivation, Institutionalization, and Productivity of Creative Neologisms*. Tübingen: Gunter Narr.
- Fragaki, Georgia, and Dionysis Goutsos. 2015. "Greek diglossia in the 20th century: A historical corpus linguistics approach". Presented at the 12th International Conference on Greek Linguistics (ICGL12), Berlin.
- Goutsos, Dionysis. 2010. "The Corpus of Greek Texts: A reference corpus for Modern Greek." *Corpora* 5 (1):29–44.
- Goutsos, Dionysis, and Georgia Fragaki. 2014. "Prosfati glossiki allagi sta ellinika: Sxediasmos tou Diaxronikou Somatos Ellinikon Kimenon tou 20^{ou} aiona." In *Selected Papers of the 11th International Conference on Greek Linguistics*, edited by G. Kotzoglou, K. Nikolou, E. Karantzola, K. Frantzi, I. Galantomos, M. Georgalidou, V. Kourti-Kazoullis, Ch. Papadopoulou, and E. Vlachou, 318–29. Rhodes: University of the Aegean.
- Hatzigeorgiu, Nikolaos, Spiliotopoulou, Sophia, Vakalopoulou, Anna, Papakostopoulou, Anna, Piperidis, Stelios, Gavriilidou, Maria, and George Carayannis. 2001. "Ethnikos Thisavros Ellinikon Kimenon (ETHEG): Soma kimenon tis Neas Ellinikis sto Diadiktyo." *Studies in Greek Linguistics* 21:812–21.
- Hunston, Susan. 2002. *Corpora in Applied Linguistics*. Cambridge: Cambridge University Press.

- Karaoglou, Ch. L. 2005. "Ta idika periodika kai i taksinomisi tous". In *La presse Grecque de 1784 à nos jours. Approches historiques et théoriques. Actes du Colloque International, Athènes, 23-25 mai 2002*, edited by Loukia Droulia, 263–75. Athens: Institute of Modern Greek Studies.
- Leech, Geoffrey, Marianne Hundt, Christian Mair, and Nicholas Smith. 2009. *Change in Contemporary English: A Grammatical Study*. Cambridge: Cambridge University Press.
- Mair, Christian. 2009. "Corpora and the study of recent change in language". In *Corpus Linguistics. An International Handbook*. Vol. 2., edited by Anke Lüdeling, and Merja Kytö, 1109–25. Berlin/New York: Walter de Gruyter.
- Marchi, Anna. 2010. "'The moral in the story': A diachronic investigation of lexicalised morality in the UK press." *Corpora* 5 (2):161–89.
- Meyer, Charles F. 2002. *English Corpus Linguistics: An Introduction*. Cambridge: Cambridge University Press.
- Nevalainen, Terttu, and Helena Raumolin-Brunberg. 2003. *Historical Sociolinguistics: Language Change in Tudor and Stuart England*. London: Routledge.
- Onelli, C. Proietti, D. Seidenari C., and F. Tamburini. 2006. The DiaCORIS project: A diachronic corpus of written Italian." In *Proceedings of the 5th International Conference on Language Resources and Evaluation, LREC 2006, Genoa*. Available at: http://hnk.ffzg.hr/bibl/lrec2006/pdf/611_pdf.pdf.
- Partington, Alan. 2010. "Modern Diachronic Corpus-Assisted Discourse Studies (MD-CADS) on UK newspapers: An overview of the project." *Corpora* 5 (2):83–108.
- Partington, Alan. 2012. "The changing discourses on antisemitism in the UK press from 1993 to 2009: A modern diachronic corpus-assisted discourse study." *Journal of Language and Politics* 11 (1):51–76.
- Schneider, Edgar W. 2004. "Investigating variation and change in written documents." In *The Handbook of Language Variation and Change*, edited by J. K. Chambers, Peter Trudgill, and Natalie, Schilling-Estes, 67–96. Oxford: Blackwell.
- Taavitsainen, Irma, Kytö, Merja, Claridge, Claudia, and Jeremy Smith 2015. *Developments in English: Expanding Electronic Evidence*. Cambridge: Cambridge University Press.

