



FACHLITERATUR
EDITION ROMIOSINI
ΕΠΙΣΤΗΜΗ



ICGL12 | 12th INTERNATIONAL CONFERENCE
ON GREEK LINGUISTICS
16 – 19 SEPTEMBER 2015
FREIE UNIVERSITÄT BERLIN, CEMOG

Proceedings of the ICGL12

vol. 2

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. 2

© 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 978-3-946142-35-5
Printed in Germany

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

ΠΕΡΙΕΧΟΜΕΝΑ

Σημείωμα εκδοτών	7
Περιεχόμενα.....	9
Peter Mackridge:	
<i>Some literary representations of spoken Greek before nationalism(1750-1801)</i>	17
Μαρία Σηφιανού:	
<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 Bacskai-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 Fragkakis, 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: Veriäten - 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 Περγέα: <i>ΣΕΠΙAME2: Μια καινούρια πηγή αναφοράς για την Ελληνική ως Γ2</i>	419
Μαρία Ιακώβου & Θωμαΐς 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 NEK 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 Kyrioti: <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 & Δώρας 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 mimos</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 Gemimates</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 Armostis: <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:	
<i>Η μεθοδολογική προσέγγιση της εικόνας της Γερμανίας στις ελληνικές εφημερίδες</i>	1095
Anastasia Tzilinis:	
<i>Begründendes Handeln im neugriechischen Wissenschaftlichen Artikel: Die Situierung des eigenen Beitrags im Forschungszusammenhang.....</i>	1109
Κυριακούλα Τζωρτζάτου, Αργύρης Αρχάκης, Άννα Ιορδανίδου & Γιώργος Ι. Ευδόπουλος:	
<i>Στάσεις απέναντι στην ορθογραφία της Κοινής Νέας Ελληνικής: Ζητήματα ερευνητικού σχεδιασμού</i>	1123
Nicole Vassalou, Dimitris Papazachariou & Mark Janse:	
<i>The Vowel System of Mišótika Cappadocian</i>	1139
Marina Vassiliou, Angelos Georganas, Prokopis Prokopidis & Haris Papageorgiou:	
<i>Co-referring or not co-referring? Answer the question!.....</i>	1155
Jeroen Vis:	
<i>The acquisition of Ancient Greek vocabulary.....</i>	1171
Christos Vlachos:	
<i>Mod(aliti)es of lifting wh-questions.....</i>	1187
Ευαγγελία Βλάχου & Κατερίνα Φραντζή:	
<i>Μελέτη της χρήσης των ποσοδεικτών λίγο-λιγάκι σε κείμενα πολιτικού λόγου</i>	1201
Madeleine Voga:	
<i>Τι μας διδάσκουν τα ρήματα της ΝΕ σχετικά με την επεξεργασία της μορφολογίας.....</i>	1213
Werner Voigt:	
<i>«Σεληνάκι μου λαμπρό, φέγγε μου να περπατώ ...» oder: warum es in dem bekannten Lied nicht so, sondern eben φεγγαράκι heißt und ngr. φεγγάρι</i>	1227
Μαρία Βραχιονίδου:	
<i>Υποκοριστικά επιρρήματα σε νεοελληνικές διαλέκτους και ιδιώματα</i>	1241
Jeroen van de Weijer & Marina Tzakosta:	
<i>The Status of *Complex in Greek.....</i>	1259
Theodoros Xioufis:	
<i>The pattern of the metaphor within metonymy in the figurative language of romantic love in modern Greek.....</i>	1275

CO-REFERRING OR NOT
CO-REFERRING?
ANSWER THE QUESTION!

Marina Vassiliou¹, Angelos Georgaras², Prokopis Prokopidis¹ & Haris Papageorgiou¹

¹Institute for Language and Speech Processing, ²University of Athens

{mvas, prokopis, xaris}@ilsp.gr, angelosgeorgaras@hotmail.com

Περίληψη

Επίλυση συναναφοράς είναι ο εντοπισμός εκφράσεων που αναφέρονται στην ίδια οντότητα. Βασική πτυχή στην Επεξεργασία Φυσικής Γλώσσας, συμβάλλει σε εφαρμογές όπως εξόρυξη πληροφορίας, περίληψη κειμένου και συστήματα ερωταποκρίσεων. Στο παρόν άρθρο περιγράφεται ένα σύστημα επίλυσης συναναφοράς, το πρώτο, εξ όσων γνωρίζουμε, για την Ελληνική. Χειρίζεται όλες τις εκφράσεις (ονοματικές και αντωνυμικές, κενά υποκείμενα) και βασίζεται σε ένα κανονιστικό αλγόριθμο, ο οποίος, εφαρμόζοντας μία σειρά – κυρίως μορφολογικών και συντακτικών – κριτηρίων δημιουργεί συστάδες συναναφερόμενων εκφράσεων. Η πρώτη υλοποίηση του συστήματος επιτυγχάνει πολύ ικανοποιητικά ποσοστά ακρίβειας (~83%). Επίσης αναπτύχθηκε σχήμα επισημείωσης σχέσεων συναναφοράς που εφαρμόστηκε σε ένα σώμα κειμένων μεγέθους >90.000 λέξεων.

Keywords: coreference, near identity, bridging, coreference resolution, coreference annotation, Natural Language Processing

1. Introduction

Coreference Resolution (CR) denotes the detection of all linguistic expressions in a text which refer to the same discourse entity. It has been (and still is) one of the most intriguing tasks within Natural Language Processing over the past twenty years (see

Ng (2010) for a concise review of the field). An interesting problem per se, coreference resolution also assumes an ancillary role in other NLP tasks such as information extraction, text summarisation, machine translation and question answering.

Coreference resolution is construed as a clustering problem, i.e. the coreferent expressions within a text must be linked as members of the same chain and grouped into the same cluster. The solution to this clustering problem has been attempted by various systems, which are either based on machine learning algorithms (Soon et al., 2001, Ng and Cardie, 2002, Versley et al., 2008b, Durrett and Klein, 2013) or follow rule-based approaches (Lappin and Leass, 1994, Haghighi and Klein, 2009, Lee et al., 2013, O'Connor and Heilman, 2013). Machine learning systems, supervised and unsupervised ones, have taken the lead over the last decade due to the availability of manually annotated corpora; yet rule-based systems still outperform them in most cases.

Models for handling coreference could be classified into three types: First, the mention¹-pair models, where expressions in a text are examined in pairs and then local decisions are assembled together (Soon et al., 2001, Ng and Cardie 2002a, Versley et al., 2008b). Alternatively, the entity-mention models highlight the importance of simultaneously checking coreference over all mentions in text which refer to an entity (Luo et al., 2004, Yang et al., 2008a, 2004, Haghighi and Klein, 2010, Lee et al., 2011). Finally, ranking models arrange all candidate antecedents of an expression based on probability (Iida et al., 2003, Yang et al., 2008b).

CR systems utilise a wide range of features, varying between simple string matching and grammatical features (part of speech, gender, number etc.) to syntactic and semantic features (e.g. agreement, binding constraints, syntactic role, selectional restrictions, θ -roles, semantic relations etc.). Furthermore, in order to enhance precision, there have been attempts to improve coreference resolution by feeding systems with knowledge of the world, the rationale being that humans also resort to this knowledge in order to process and understand texts. This knowledge is drawn from web resources such as Wikipedia, DBPedia or YAGO (cf. Uryupina 2006, Ponzetto et al., 2007, Versley et al., 2008a, Bryl et al., 2010).

In this paper we describe a coreference resolution system, the first one, to our knowledge, that handles Greek data. The system implements a rule-based algorithm, which falls under the category of the entity-mention models, and targets both precision and recall. More specifically, with the aim of maximising recall, it initially extracts all lin-

¹ *Mention* denotes the occurrence of an expression in a text.

guistic expressions from a text and creates corresponding clusters (i.e. coreference chains). Subsequently, a series of precision-oriented criteria, which primarily rely on morphological and syntactic features, are applied in a stepwise mode, checking the consistency of already formed clusters and governing the creation of new ones.

The linguistic expressions handled by the system comprise only nominal and pronominal referential expressions as well as null subjects in both finite and gerundive clauses. Adverbial expressions with a referential function are not currently supported.

Furthermore, the current article reports on related resources that have been created in parallel, namely a coreference annotation scheme for the Greek language that also addresses bridging and near identity relations, as well as a Greek corpus annotated with these relations.

The paper is structured as follows: Section 2 provides information on the coreference resources for the Greek language, namely the annotation scheme and the corresponding corpus. Section 3 describes the coreference resolution system, while Section 4 reports on the evaluation results obtained. The concluding section outlines future directions with respect to the development of the system and the improvement of its performance.

2. Coreference annotation scheme for the Greek language

2.1 Defining relations between expressions

The Greek coreference annotation scheme specifies three types of relations between linguistic expressions, (i) coreference, (ii) near identity and (iii) bridging.

Coreference denotes a relation between (at least) two linguistic expressions, which have the same referent, i.e. they point to the same discourse entity. The linking of these expressions results in a chain, either within the sentence boundaries (examples 1-3) or spanning more than one sentence (example 4):

- (1) Η εξέλιξη αυτή ανακοινώθηκε ταυτόχρονα στις ΗΠΑ και τη Βόρεια **Κορέα** και προέκυψε δύο μήνες μετά το θάνατο του ηγέτη της Βόρειας **Κορέας**.

‘This development was announced in the US and North Korea at the same time and came two months after the death of the leader of North Korea.’

- (2) Στην Ουάσιγκτον, ο Λευκός Οίκος εξέφρασε ικανοποίηση για την απόφαση της **Πιονγκγιάνγκ**, χαρακτηρίζοντάς **την** θετικό πρώτο βήμα.

'In Washington, the White House welcomed the decision of Pyongyang and characterised it as a positive first step.'

- (3) Πέρασε ο ίδιος νύχτα στο **Ακρωτήριο όπου** ύψωσε την ελληνική σημαία.

'He went himself to Akrotiri during the night and raised the Greek flag.'

- (4) Ο Αμερικανός **Πρόεδρος** δήλωσε ... ότι “προτιμά το θέμα να επιλυθεί μέσω της διπλωματίας.” Ο Μπαράκ **Ομπάμα** ανέφερε ότι «ήδη υπάρχουν πολλές αναφορές για πόλεμο με το Ιράν.”

'The US president said ... that “the issue will preferably be resolved through diplomacy.” Barack Obama said that “there are already many references to war with Iran”’

Near identity denotes a quasi coreferential relation between two different linguistic expressions or instances of the same expression, which implicitly refer to the same entity or represent an aspect of the same entity. These expressions “are partially the same in that they share most of the important characteristics, but differ in at least one crucial dimension” (Recasens and Martí, 2010: 151). In example (5) the two instances of the expression “**ΗΠΑ**” are quasi coreferent, since the first one has a locative reading, whereas the second one refers to a related aspect of this location, namely the US government.

- (5) Η νίκη του καλωσορίστηκε τόσο στην Κούβα όσο και στις **ΗΠΑ**. ... Έτσι οι **ΗΠΑ** έγιναν πολύ εχθρικές έναντι του Κάστρο κατά το 1959.

'His victory was welcomed in Cuba as well as in the United States ... Therefore, in 1959 the United States became pretty hostile towards Castro.'

Bridging denotes an associative relation between two non-coreferent linguistic expressions. This relation may come in various types (for instance set-subset, part-whole, entity-property, contrastive association etc.). Example (6) instantiates a set-subset relation holding between the two bold-faced expressions.

(6) Σε όλες σχεδόν τις κατοικημένες περιοχές του κέντρου υπάρχει η ίδια εικόνα ερήμωσης.

‘One sees the same picture of desolation in almost all populated downtown areas.’

2.2 *The annotation scheme in a nutshell*

The annotation scheme is largely based on the PDT 2.0 (Prague Dependency Treebank)² scheme, used for manual annotation of texts at multiple levels (morphology, syntax, semantics and pragmatics); yet it contains several modifications pertinent to the Greek language, while it was also extended to accommodate near identity relations.

All kinds of expressions are taken into account (nominal expressions and pronominal elements, anaphoric adverbs, and null subjects). Notably, as regards coreference, split antecedents are allowed. Additionally, metonymy, which functions as a coreference mechanism, is also covered in the annotation scheme. Furthermore, near identity relations come in five types. Finally, nine bridging relation types are specified.

A comprehensive description of the scheme together with annotation guidelines is available online³.

2.3 *The annotated corpus*

The resource annotated with the aforementioned scheme is a medium-sized corpus comprising mainly news texts drawn from various sources. Table 1 on the next page illustrates the corpus profile.

Figure 1 illustrates the TrEd⁴ annotation tool that was used for the annotation of the corpus. Text is represented as sequences of dependency trees and the coreferent nodes are linked with an arrow.

2 <https://ufal.mff.cuni.cz/pdt2.0/>

3 <http://gdt.ilsp.gr/guidelines/coreference/coreferenceannotationguidelines.pdf/@@download/file/coreferenceannotationguidelines.pdf>

4 <http://ufal.mff.cuni.cz/tred/>

Source	Texts	Sentences	Tokens
newspapers	11	186	4.120
elwikinews	92	675	14.739
wikipedia	11	263	6.263
europarl	45	1.301	30.076
setimes	6	67	1.684
voa	9	89	2.551
in.gr	81	791	21.843
opinion texts	20	591	11.788
Total	275	3.963	93.064

Table 1 | Coreference-annotated corpus

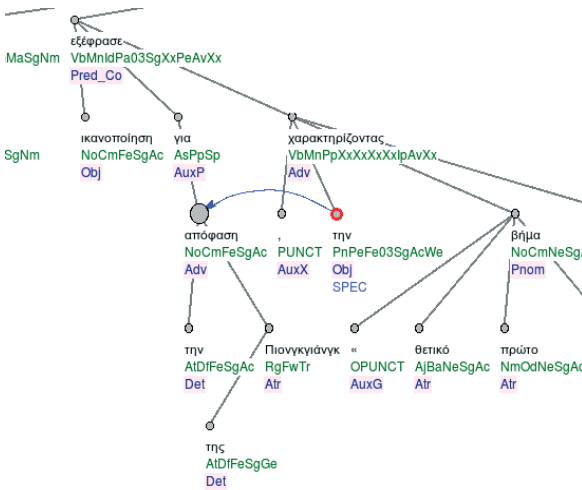


Figure 1 | Screenshot of the annotation environment

3. A coreference resolution system for Greek

Our coreference resolution system for Greek belongs to the rule-based paradigm rather than the machine learning one. This decision was dictated not only by the scarcity of available training resources in Greek, but rather by the fact that rule-based approaches yield the best results as yet. Following the entity-mention model of dealing with coreference, the system attempts to jointly detect all the expressions in a text, nominal and

pronominal ones, phonologically instantiated or not, which refer to the same entity.

The system receives as input an XML file (see Figure 2) containing text annotated with information about token, lemma, Part-of-Speech (PoS) tag and syntactic dependency relations.

```

<children id="t-20140701-ingr-8013-s1w7">
<t_lemma>Ρωσία</t_lemma>
<deepord>7</deepord>
  <children id="t-20140701-ingr-8013-s1w6">
    <t_lemma>ο</t_lemma>
    <deepord>6</deepord>
    <form>της</form>
    <tag>AtDfFeSgGe</tag>
    <sentord>6</sentord>
    <afun>Det</afun>
  </children>
<form>Ρωσίας</form>
<tag>NoPrFeSgGe</tag>
<sentord>7</sentord>
<afun>Atr</afun>
</children>

```

Figure 2 | Indicative input to the CR system

s1w6	ο	της	AtDfFeSgGe	Det	
s1w7	Ρωσία	Ρωσίας	NoPrFe03SgGe	Atr	(7)
s1w8	λόγω	λόγω	AsPpSp	AuxP	
s1w9	ο	των	AtDfFePlGe	Det	
s1w10	κύρωση	κυρώσεων	NoCmFe03PlGe	Atr	(4)
s1w11	.	.	PTERM_P	AuxK	
s2w1	ο	Οι	AtDfFePINm	Det	
s2w2	κύρωση	κυρώσεις	NoCmFe03PINm	Sb	(4)
s2w3	που	που	PnReFe03PINmXx	Sb	(4)
s2w4	έχω	έχουν	VbMnIdPro3PlXlpAvXx	AuxV	
s2w5	επιβάλλω	επιβληθεί	VbMnNfXxXxXxPePvXx	Atr	
s2w6	σου	στη	AsPpPaFeSgAc	AuxP	
s2w7	Ρωσία	Ρωσία	NoPrFe03SgAc	IObj	(7)
s2w23	μπορώ	μπορεί	VbIsIdPro3SgXlpAvXx	Obj_Co	
s2w24	να	να	PtSj	AuxV	
s2w25	οδηγώ	οδηγήσουν	VbMnIdXx03PlXxPeAvXx	Sb	
s2a1	#Cor	-	03PINm	Sb	(4)

Figure 3 | Indicative output of the CR system

The system output includes an extra annotation layer whereby each expression is flagged with the id of the cluster it has been classified into. In Figure 3, the expressions marked with the same id (see last column) are grouped into the same cluster, thus they are coreferent.

The CR algorithm is entity-oriented, trying to jointly discover all the expressions that refer to the same entity. Eight phases can be identified:

Phase 1 – Process text: Each word of the input text is represented as a 10-feature vector.

[s6w10, απεργία, NoCmFe03SgGe, 10, Obj, s6w9, 6, 11, 92, απεργίας]

- Feature 1: word id
- Feature 2: word lemma
- Feature 3: PoS tag
- Feature 4: position of the word in the sentence
- Feature 5: dependency relation tag
- Feature 6: parent node id
- Feature 7: sentence id
- Feature 8: clause id
- Feature 9: position of the word in the text
- Feature 10: token

Phase 2 – Identify expressions: All candidate expressions in a text are identified and classified into four categories, namely (i) nominal mentions, (ii) pronominal mentions, (iii) null subjects of a finite clause and (iv) null subjects in non-finite clauses. Expressions which lack anaphoric properties, that is, they do not refer to an entity, for example attributive expressions, are not considered eligible and are excluded.

Phase 3 – Initial clustering: All nominal mentions are clustered based on their lemma. So mentions with the same lemma are considered coreferent and grouped together. Mentions which differ in lemma yet are joined in an appositional structure, for instance “*Ο Βλαντιμίρ Πούτιν, ο πρόεδρος της Ρωσίας*”, are also considered coreferent and are clustered together.

Phase 4 – Refine clustering: The initial clustering of the nominal mentions is re-examined and modified accordingly on the basis of pre- or post- modifiers. So, mentions that have been grouped together due to lemma identity, will no longer form the same

cluster if their modifiers differ (e.g. “*αρχική έκθεση*” vs. “*τελική έκθεση*”).

Phase 5 – Handle null subjects: For each null subject the algorithm detects its coreferent mention, which could itself be a null subject, utilising mainly agreement and structural features. If the mention already belongs to a cluster, then the null subject is included too. Otherwise, a two-member cluster is formed comprising the null subject and its antecedent.

Phase 6 – Handle pronominal mentions: The algorithm detects the antecedent of each pronominal element utilising agreement and locality features. For instance, the antecedent of a relative pronoun must precede the pronoun and agree with it in gender and number. In a similar vein to Phase 5, if the antecedent has been already clustered, then it pied-pipes the pronominal mention to its own cluster. In a different case, the pronoun and its antecedent form a two-member cluster.

Phase 7 – Merge clusters: The algorithm re-examines the already-formed clusters and merges the ones that contain coreferent items. For example, the two-member cluster {*Πούτιν, Πρόεδρος*} will be merged with a cluster containing mentions which have the lemma “*πρόεδρος*”.

Phase 8 – Handle residual items: Further merging is attempted concerning smaller clusters which contain null subjects and pronouns. Additionally, singletons, i.e. one-member clusters, are removed.

4. Evaluation of the coreference resolution system

4.1 Evaluation setup

The evaluation of the system was performed on a corpus of texts drawn from the web and the Greek Dependency Treebank⁵. Two subsets (Sets 1 and 2), constituting the largest part of the corpus, belonged to the news domain, while a small corpus subset (Set 3) comprised opinionated texts. The corpus was automatically annotated with information about lemma, PoS tag and dependency relations. The annotations were then manually corrected so as to establish that the CR system’s performance could be assessed without being affected by any pre-processing errors, while the remaining errors could be attributed to the algorithm itself.

5 <http://gdt.ilsp.gr>

For evaluation purposes it is required to compile a golden corpus, i.e. the correct output of the system, against which the system output is compared. To obtain the golden output, a second version of the corpus was created, where the coreferential chains were manually indicated.

Table 2 provides the details of the evaluation corpus.

Sets	Type	Source	Texts	Sentences	Mean number of sentences	Tokens	Mean number of tokens
Set1	news	GDT; web	40	311	7,78	6.913	22,23
Set2	news	in.gr	81	791	9,77	21.843	27,61
Set3	restaurant reviews	athinorama	15	446	29,73	8.711	19,53
			136	1.548		37.467	

Table 2 | Evaluation corpus

4.2 Evaluation results

A CR system's performance is assessed in terms of its capacity to detect the coreferent expressions within a text and cluster them together. For our experiments a series of established evaluation metrics have been employed, namely MUC (Vilain et al., 1995), B³ (Bagga and Baldwin, 1998), CEAf (Luo, 2005) and BLANC (Recasens and Hovy, 2011, Luo et al., 2014). These metrics calculate the system's felicity in detecting coreference relations, by comparing the system output to the golden output, and produce recall, precision and *F* scores.

Recall, in this case, measures the number of the correct coreference chains (clusters) that the system succeeds in creating. Precision, on the other hand, denotes how many of the system-created coreference chains (clusters) are correct, while the *F* score indicates the harmonic mean of these two measures. It should be noted that at this point, where the first system version was evaluated, it was necessary to establish that the system was capable of creating correct clusters; subsequently, the precision score was more important than recall.

Table 3 lists the results obtained for the specific dataset. The bold-faced values are the highest ones obtained per set regarding the precision and *F* scores. It is observed that

the average precision for all sets approximates 83%, while the corresponding *F* score is close to 77%.

	Set1			Set2			Set3		
Metrics	Recall	Precision	<i>F</i>	Recall	Precision	<i>F</i>	Recall	Precision	<i>F</i>
MUC	0,78	0,92	0,84	0,66	0,83	0,73	0,52	0,72	0,60
B ³	0,74	0,92	0,81	0,61	0,82	0,69	0,52	0,73	0,60
CEAFM	0,79	0,88	0,83	0,68	0,78	0,73	0,59	0,72	0,65
CEAFE	0,83	0,82	0,82	0,75	0,71	0,72	0,68	0,66	0,67
BLANC	0,69	0,91	0,77	0,53	0,78	0,61	0,36	0,66	0,45

Table 3 | Evaluation results for the first release of the CR system

From the evaluation results it is evident that the first system release has achieved a promising performance, despite the fact that some phenomena, for instance split antecedents, metonymy or long-distance or inter-sentential dependencies, have by design not been addressed in the first implementation. Therefore, it is expected that there is ground for further improvement.

Furthermore, some preliminary testing has shown that factors such as text size or number of mentions possibly have an impact on the system's performance (see Table 4); yet further experimentation with more texts and various domains is needed to establish this correlation and its statistical significance. In a similar vein, the high frequency-of-occurrence of null subjects, which is observed in Set 3, predictably influences the system's precision, as the task of finding the antecedent of null subjects is more intricate.

	Set1	Set2	Set3
Precision	0,92	0,83	0,73
<i>F</i>	0,84	0,73	0,67
Mean number of sentences	7,78	9,77	29,73
Mean number of tokens	22,23	27,61	19,53
Mean number of mentions	0,40	0,18	0,13
Mean number of null subjects	0,10	0,12	0,24

Table 4 | System's precision in relation to the evaluation corpus size

5. Summary and future directions

In the present paper a coreference resolution system has been presented. To our knowledge, it is the first system that has been developed for the Greek language. It follows a domain-independent, rule-based approach that falls under the category of the entity-mention models. The system can handle all types of expressions, nominal and pronominal ones as well as null subjects.

High precision being mainly targeted, clustering of the coreferent expressions in a text is performed gradually through the employment of a set of morphological and syntactic criteria. The first system release, reported on here, has yielded promising results, having achieved precision of ~83%.

In parallel with the CR system, an annotation scheme has been developed, which apart from coreference, describes two types of relations between linguistic expressions: near identity and bridging. The scheme has been used for annotating a medium-size corpus of Greek texts mainly originating from the news domain.

Future directions involve handling of more demanding phenomena such as split antecedents or the detection of inter-sentential pronoun-antecedent chains. The anaphoricity check, through which expressions not referring to an entity are excluded from being possible candidates, could also be enhanced. Furthermore, exploiting semantic knowledge and/or the factor of context (cf. Recasens et al., 2013: 2) is expected to substantially improve the system's precision in detecting coreferential expressions. Finally, it is planned to experiment with different genres and domains in order to assess whether the domain/genre factor has an impact on performance, even though the system's rationale and implementation are domain-independent.

Acknowledgement

The research leading to the results above has been funded by the «POLYTROPON» project (ΟΠΣ: 448306, ΠΔΕ: 2013ΣΕΕ01380028).

References

- Bagga, Amit, and Breck Baldwin. 1998. "Algorithms for Scoring Coreference Chains". In *Proceedings of the 1st International Conference on Language Resources and Evaluation (LREC'98), Workshop on Linguistic Coreference*, Granada, Spain, 563–566.
- Bryl, Volha, Claudio Giuliano, Luciano Serafini, and Kateryna Tymoshenko. 2010. "Supporting Natural Language Processing with Background Knowledge: Coreference Resolution Case". In *Proceedings of the 9th International Semantic Web Conference (ISWC 2010)*, Shanghai, China, 80–95.
- Durrett, Greg, and Dan Klein. 2013. "Easy Victories and Uphill Battles in Coreference Resolution". In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing (EMNLP 2013)*, Seattle, USA, 1971–1982.
- Haghighi, Aria, and Dan Klein. 2009. "Simple Coreference Resolution with Rich Syntactic and Semantic Features". In *Proceedings of the EMNLP 2009, Conference on Empirical Methods in Natural Language Processing*, Singapore, 1152–1161.
- Haghighi, Aria, and Dan Klein. 2010. "Coreference Resolution in a Modular, Entity-Centered Model". In *Proceedings of Human Language Technologies: The 11th Annual Conference of the North American Chapter of the Association for Computational Linguistics (HLT-NAACL 2010)*, Los Angeles, CA., USA, 385–393.
- Iida, Ryu, Inui Kentaro, Takamura Hiroya, and Yuji Matsumoto. 2003. "Incorporating Contextual Cues in Trainable Models for Coreference Resolution". In *Proceedings of "The Computational Treatment of Anaphora" workshop held in conjunction with the 10th Conference of the European Chapter of the Association for Computational Linguistics (EACL 2003)*, Budapest, Hungary, 23–30.
- Lappin, Shalom, and Herbert J. Leas. 1994. "An Algorithm for Pronominal Anaphora Resolution". *Computational Linguistics* 20(4): 535–561.
- Lee, Heeyoung, Angel Chang, Yves Peirsman, Nathanael Chambers, Mihai Surdeanu, and Dan Jurafsky. 2013. "Deterministic Coreference Resolution Based on Entity-Centric, Precision-Ranked Rules". *Computational Linguistics* 39(4): 885–916.
- Luo, Xiaoqiang, Abe Ittycheriah, Hongyan Jing, Nanda Kambhatla, and Salim Roukos. 2004. "A Mention-Synchronous Coreference Resolution Algorithm Based on the Bell Tree". In *Proceedings of the 42nd*

Annual Meeting of the Association for Computational Linguistics, Barcelona, Spain, 135–142.

- Luo, Xiaoqiang. 2005. “On Coreference Resolution Performance Metrics”. In *Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing (HLT '05)*, Vancouver, British Columbia, Canada, 25–32.
- Luo, Xiaoqiang, Sameer Pradhan, Marta Recasens, and Eduard Hovy. 2014. “An Extension of BLANC to System Mentions”. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, Baltimore, Maryland, USA, 24–29.
- Ng, Vincent, and Claire Cardie. 2002. “Improving Machine Learning Approaches to Coreference Resolution”. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL-02)*, Philadelphia, USA, 104–111.
- Ng, Vincent. 2010. “Supervised Noun Phrase Coreference Research: The First Fifteen Years”. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics (ACL '10)*, Uppsala, Sweden, 1396–1411.
- O'Connor, Brendan, and Michael Heilman. 2013. “ARKref: A Rule-Based Coreference Resolution System”. ArXiv: 1310.1975. <https://arxiv.org/abs/1310.1975>
- Ponzetto, Simone Paolo, and Michael Strube. 2007. “Knowledge Derived From Wikipedia For Computing Semantic Relatedness”. *Journal of Artificial Intelligence Research* 30:181–212.
- Recasens, Marta, and Antònia M. Martí. 2010. “AnCora-CO: Coreferentially Annotated Corpora for Spanish and Catalan”. *Language Resources and Evaluation* 44(4):315–345.
- Recasens, Marta, and Eduard Hovy. 2011. “BLANC: Implementing the Rand Index for Coreference Evaluation”. *Natural Language Engineering* 17(4):485–510.
- Recasens, Marta, Matthew Can, and Dan Jurafsky. 2013. “Same Referent, Different Words: Unsupervised Mining of Opaque Coreferent Mentions”. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL 2013)*, Atlanta, Georgia, USA, 897–906.
- Soon, Wee Meng, Hwee Tou Ng, and Daniel Chung Yong Lim. 2001. “A Machine Learning Approach to Coreference Resolution of Noun Phrases”. *Computational Linguistics* 27(4):521–544.
- Uryupina, Olga, 2006. “Coreference Resolution with and without Linguistic Know-

- ledge”. In *Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC 2006)*, Genoa, Italy, 893–898.
- Versley, Yannick, Simone Ponzetto, Massimo Poesio, Vladimir Eidelman, Alan Jern, Jason Smith, Xiaofeng Yang, and Alessandro Moschitti. 2008a. “BART: A Modular Toolkit for Coreference Resolution.” In *Proceedings of the 6th International Conference on Language Resources and Evaluation (LREC 2008)*, Marrakech, Morocco, 962–965.
- Versley, Yannick, Alessandro Moschitti, Massimo Poesio, and Xiaofeng Yang. 2008b. “Coreference Systems based on Kernel Methods.” In *Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)*, Manchester, UK, 961–968.
- Vilain, Marc, John Burger, John Aberdeen, Dennis Connolly, and Lynette Hirschman. 1995. “A Model Theoretic Coreference Scoring Scheme”. In *Proceedings of the 6th Conference on Message Understanding (MUC6 ‘95)*, Columbia, Maryland, 45–52.
- Yang, Xiaofeng, Su Jian, Jun Lang, Tan Chew Lim, Ting Liu, and Sheng Li. 2008a. “An Entity-Mention Model for Coreference Resolution with Inductive Logic Programming. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies (ACL-HLT 2008)*, Columbus, Ohio, USA, 843–851.
- Yang, Xiaofeng, Su Jian, GuoDong Zhou, and Tan Chew Lim. 2004. “An NP-Cluster Based Approach to Coreference Resolution”. In *Proceedings of the 20th International Conference on Computational Linguistics (COLING 2004)*, Geneva, Switzerland, 226–232.
- Yang, Xiaofeng, Su Jian, and Tan Chew Lim. 2008b. “A Twin-Candidate Model for Learning-Based Anaphora Resolution”. *Computational Linguistics* 34(3):327–356.

