

COMPUTATIONAL ANALYSIS OF SUFFIXES AND BOUND STEMS OF THE GREEK LANGUAGE: A CRASH TEST WITH *LINGUISTICA*

ATHANASIOS KARASIMOS & EVANTHIA PETROPOULOU

Abstract

This paper attempts to participate in the ongoing discussion in search of a suitable model for the computational treatment of Greek morphology. Focusing on the unsupervised morphology learning technique, and particularly on the model of *Linguistica* by Goldsmith (2001), we attempt a computational treatment of specific word formation phenomena in Modern Greek (MG), such as suffixation and compounding with bound stems, through the use of various corpora. The inability of the system to receive any morphological rule as input, hence the term 'unsupervised', interferes to a great extent with its efficiency in parsing, especially in languages with rich morphology, such as MG, among others. Specifically, neither the rich allomorphy, nor the complex combinability of morphemes in MG appear to be treated efficiently through this technique, resulting in low scores of proper word segmentation (22% in inflectional suffixes and 13% in derivational ones), as well as the recognition of false morphemes.

Key words: Unsupervised Morphology Learning, Goldsmith, *Linguistica*, Greek derivational affixes, Greek bound stems

1. A Brief Introduction to Computational Morphology

Computational morphological analysis has already passed the half century mark, as the first attempts saw the light of day already within the earliest work on Machine Translation. Roark and Sproat (2007) honor the work of Andron (1962), Woyna (1962), Bernard-Georges et al. (1962), Boussard and Berthaud (1965), Vauquois (1965), Schweiger and Mathe (1965), Matthews (1966), Brand et al. (1969) and Hutchins (2001), among others. Over these years, many applications have been implemented; including a variety of stemmers, parsers, spelling correctors, text input systems and natural language generation systems.

Despite the importance of the afore-mentioned pioneers, the most interesting and audacious work in computational linguistics was the approach depending on finite state methods. The most dominant finite-state morphology has been the approach of Koskenniemi (1983), based on finite state transducers. Koskenniemi (1983) has implemented the theoretical investigation of Kaplan and Kay of Xerox PARC₂ in KIMMO, a two-level morphological analyzer, which still remains a state-of-the-art application in computational linguistics. Alternative approaches to computational morphology are either based on explicitly finite-state models of morphotactics (Allen et al. 1987, Roark and Sproat 2007) or models involving “suffix stripping” (McIlroy 1982, Porter 1980).

As Roark and Sproat (2007: 102) point out, Koskenniemi's Two-Level

Morphology approach has been used to develop morphological analyzers for several languages, including Modern Greek (see Σγάμπας 1997, Μαρκόπουλος 1998, Ralli & Galiotou (2001, 2003, 2005). KIMMO contains three central elements: a) *trees*, the representations of dictionaries, b) *continuation lexica*, the representation of morphological concatenation, and c) *finite-state transducers*, which implement the surface–lexical morphophonological correspondences, changes and processes.

2. Unsupervised Morphology Learning: A theoretical approach

2.1. An Introduction to Unsupervised Morphology Learning

As opposed to the computational analyses on syntax, computational work on morphology has been relatively scarce. According to Roark and Sproat (2007), the absence of a corpus of morphologically annotated words put a burden on the development of a machine learning morphological system that could rival a morphologically–complex analyzer such as the one proposed by Koskenniemi (1983). However, close to the dawn of the new millennium, the interest in statistical models of morphology, particularly of unsupervised (or lightly supervised) morphology–learning from annotated corpora, has rapidly increased. Special attention has been paid to automatic – basically unsupervised – methods for the discovery of morphological alternations.

In order to give a clear picture of this system, it would be best to provide a definition of what morphological learning is, through a number of illustrative examples. One of the main objectives of the system is to discover relationships between words sharing common strings, on the basis of specific data. Take for instance, the word *άνθρωπος* (*anthropos*) (‘man’) and its alternative inflectional forms, e.g. *ανθρώπου* (*anthropu*), *άνθρωπο* (*anthropo*), *άνθρωπε* (*anthrope*), *άνθρωπι* (*anthropi*), *ανθρώπων* (*anthropon*) and *ανθρώπους* (*anthropus*). The system tries to create a set of words with related forms. Another aim is to generate words on the basis of some regular suffixation (inflectional or derivational) pattern; the noun *μωρό* (*moro*) (‘baby’ diminutive) would become *μωράκι* (*moraki*) (‘little baby’) by adding the derivational suffix *–άκι* (*–aki*). The goal is to derive new morphological forms never encountered before, via the application of a set of rules (Roark and Sproat 2007). However, allomorphy poses a serious problem for both tasks. By treating allomorphy, the goal is to find related morphological forms of the same word, such as *κύμα* and *κύματα* (*kima~ kimat(a)*) (‘wave’), which are not the product of any phonological and morphological rules.

Since most of the recent research has been carried out within the field of unsupervised morphological learning, we will focus our discussion and criticism on this system, and specifically on the theory of *Minimum Length Description* (MLD) proposed by Goldsmith (2001) [other recent works in the same direction are Yarowsky and Wicentowski 2001, Schone and Jurafsky 2001, Creutz and Lagus 2002]. Goldsmith’s (2001) theory and the implementation of his program *Linguistica*

are based on the framework of Rissanen's (1989) MDL. His paper is not the first work on unsupervised morphology learning, as there are three other approaches by previous researchers. However, his work is certainly the mostly cited, and is considered to be the standard model compared to other systems.

Research in automatic morphological analysis has been divided into four major approaches. The first approach by Harris (1955, 1967), Hafer and Weiss (1974) tries to identify morpheme boundaries and then classify them into stems, prefixes and suffixes. They attempt to use prefix/ suffix conditional entropy in order to set boundaries inside words. The second approach proposes bigrams and trigrams as parts of a morpheme's form (see Janssen 1992, Klenk 1992, Flenner 1994, 1995). According to this hypothesis, the local information, i.e. the summary of probability and frequencies, in a string of two or three phonemes is sufficient to set boundaries. The third approach focuses on the pattern of phonological relationships between pairs of related words, as shown by Dzeroski and Erjavec 1997. Their goal is to predict the form of a word based on morphological principles and a given word form. Finally, the fourth approach by Goldsmith (2001) will be discussed in detail in the next section.

2.2. Goldsmith's *Minimum Length Description* (2001)

Goldsmith's system starts with a very large corpus of annotated texts and produces a range of *signatures* along with words that belong to these signatures. A *Signature* is a set of affixes (prefixes or suffixes) that combine with a given set of stems (Goldsmith 2001, Roark and Sproat 2007). An example suffix signature in English could be *NULL.ed.ing.s*, which combines with the stems *jump*, *laugh*, *walk*, *talk*, etc., all of which take the signature's suffixes in order to create words, such as *jumpø*, *jumped*, *jumping* and *jumps*. Other examples of signatures are *e.ed.ing*, *NULL.s*, *NULL.ing.s*, *NULL.er.est.ly*, etc. As we can see, these signatures are like paradigms, but they usually contain both inflectional and derivational suffixes. So the basic schema of how signatures work is the following:

$$(1) \left\{ \begin{array}{l} \text{STEM}_1 \\ \text{STEM}_2 \\ \text{STEM}_3 \end{array} \right\} \left\{ \begin{array}{l} \text{SUFFIX}_1 \\ \text{SUFFIX}_2 \end{array} \right\}$$

A closer look at the signatures reveals that the sets are not always complete. Usually the past tense suffixes are absent, even for regular verb stems. For example, Roark and Sproat (2007:120) point out that the signature *NULL.er.ing.s* proposed by Goldsmith (2001: 179), which includes stems such as *blow*, *broadcast*, *drink*, *feel* does not display the *-ed* suffix, since the verbs are irregular in their past tense form. However, the *-ed* suffix is also absent from stems such as *bomb* and *farm*, which, although are regular in their past tense form (*bombed* and *farmed*), (but which) did not unfortunately occur in the corpus! Goldsmith discusses in general terms some

problems with signatures and notes that his system is incapable of handling alternations (e.g. *allomorphs*), such as *feel/felt*, since it deals only with affixation.

As it will be demonstrated in the next section, this kind of allomorphic alternation can be an enormous problem, if one tries to apply an Unsupervised Morphology Learning Model (UMLM) for example to the Greek language, which exhibits a high degree of complex allomorphy in every word formation process (inflection, derivation, compounding). The combinability of derivational suffixes and bound stems deteriorates the problem even more.

2.2.1 Candidate generation and Evaluation

The creation of signatures involves two steps: first, the system generates a number of candidate signatures (assigning them a membership) and then evaluates the candidates. For candidate generation, the segmentation method is based on *weighted mutual information*. This method starts creating a list of affixes, an inverse lexicon (starting from the right edge of words), and builds a set of possible suffixes up to the length of six phonemes (for example –ούτσικ[ος]/ μικρ#ουτσικ#ος (mikrutsikos) ‘very small’). It then uses an algorithm that weighs all the possible suffixes in order to obtain real suffixes, and groups them into a signature. Here, Goldsmith proposes an evaluation metric based on minimum length description, whereby the best proposal for the signatures is the one which includes the most compact description of the corpus/ language.

2.2.2 Criticism

As Roark and Sproat (2007:123) correctly point out, Goldsmith’s method is “*the de facto gold standard for work on unsupervised acquisition of morphology*”. However, this system is still a far cry from perfection. As already observed, an UMLM does not use morphological and phonological rules, does not have a pre-built lexicon, and obviously does not take advantage of any linguistic (more specifically morphological) theory or framework. It only tries to split words on the basis of huge corpora. Several researchers complain that Goldsmith’s method does not exploit semantic and syntactic information. This criticism echoes the psycholinguistic approach and its objection to the fact that children and adults access other information besides the set of stems and affixes. However, considering the fact that even morphological rules or theories are left out of the model, it would perhaps be too much to anticipate the use of semantic and structural information.

The failure to correctly segment words into actual morphemes is due to the lack of morphological and phonological rules, the non-use of Lexical Phonology and the occurrence of rare, marked and irregular cases. This can happen on both the orthographical and phonological levels of word transcription:

- | | | |
|-----|----------------------------------|--|
| (2) | έγραψα >
'I wrote'
εγραψα> | ε – γραφ – σ(α) [dissimilation]
stem: γραφ
ε – γραφ – s(a) |
|-----|----------------------------------|--|

‘I wrote’ stem: γραφ

Parsing failure is more frequent in morphologically rich languages, such as Greek, Finnish, Swedish, Hungarian and Turkish. The high productivity of compounding and derivation complicates things more, introducing the factor of affix combinability. According to Kurimo *et al* (2007), the highest score of an UML model evaluation for Finnish and Turkish was 65% and 64% respectively, and the lowest score was 3% and 2%, in spite of the fact that Kurimo’s system was partly assisted by supervised morphology. One would expect that the application of the model to Greek would result in an even lower score, due to the extensive allomorphy of the language (see Karasimos 2001, Ralli 2005, 2007), as well as the complex combinatorial properties of affixes and bound stems. Melissaropoulou (2007a, 2007b) and Melissaropoulou & Ralli (2008) note that in Greek, a sequence of as many as five derivational suffixes in a row may be found within the same word.

(3)	χορευταρούλικο	>	χορ – ευ – τ – αρ – ουλ – ικ – (ο)
	(xoreftaruliko)		
	‘little great dancer’		stem – ds ¹ – ds – ds – ds – ds – (is)
	κοινωνικότητα	>	κοιν – ων – ι – ικ – οτητα(ø)
	(kinonikotita)		
	‘sociability’		stem – ds – ds – ds – ds – (is)
	ποτιστικός	>	ποτ – ισ – τ – ικ (ος)
	(potistikos)		
	‘watering (adj)’		stem – ds – ds – ds – (is)
	ξαναεπαναλαμβάνω	>	ξανά – επανα – λαμβάν(ω)
	(ksanaepanalamvano)		
	‘repeat again’		dp – dp stem
	παρασυμπαραστέκομαι	>	πάρá – συν – παρα – στεκ(ομαι)
	(parasimbarastekome)		
	‘aid (sb) too much’		dp – dp – dp stem
	συμπεριφέρομαι	>	συν – περί – φέρ(ομαι)
	(simberiferome)		
	‘behave’		dp – dp stem

Going back to Goldsmith’s theory, a signature is a set of suffixes that can be attached to a set of stems. Therefore, one should create signatures of suffixes that combine with other signatures. It is easy to imagine how complex a system with a net of suffix/ prefix signatures can become; the selection restrictions and combinational choices of derivational suffixes and bound stems render the creation of these signatures almost impossible or completely defective.

¹ DP = derivational prefix, DS = derivational suffix, IS = inflectional suffix

3. Derivational suffixes vs. Bound stems

3.1. Allomorphy and short overview of previous work

As already pointed out, allomorphy is a serious problem for UML models and an issue that almost no one in computational morphology tries to solve or even discuss. Allomorphs are different forms of the same morpheme that share lexical information, but differ unpredictably and arbitrarily in their phonological form and in the morphological environment, where they appear. Allomorphy is a central issue in morphology; however apart from a few exceptions it has never become the focus of attention, particularly within the generative grammar framework. As Ralli (2006: 2) claims “*the reason for such neglect is mainly the fact that allomorphy is usually considered as nothing more than the absence of uniformity, resulting either from historical processes or from borrowing*”.

Lieber (1982), Carstairs (1987), Booij (1997), and Ralli (1994, 2000, 2005, 2006) provide a thorough treatment of allomorphy proposing various analyses and raising several interesting points; their approaches deal with the problem from a morphological point of view. In particular, Ralli shows that the systematic allomorphic behavior of a number of Greek stems affects the organization of paradigms in a significant manner. On the contrary, Mascaro (1996, 2007), Thornton (1997), Galani (2003) and Drachman (2006) analyze allomorphy on the basis of phonological theories. Moreover, Karasimos (2001) provides a wide range of examples in all three word-formation processes, inflection, derivation and compounding, and shows how important allomorphy can be in the Greek language.

3.2. Derivational prefixes and suffixes

Affixes, depending on their position with respect to a stem/root, are distinguished into prefixes and suffixes. The prefixes are a small group of morphemes, the majority of which used to belong to the class of prepositions of Ancient Greek; some of them still participate in lexicalized phrases, such as *ανά έτος* (ana etos) ‘per annum’, *συν τοις άλλοις* (sin tis alis) ‘moreover’. Only 32% of the prefixes display allomorphic behaviour. This allomorphy is mostly due to certain phonological rules that became inactive in Modern Greek, such as Grassman’s Law or the aspiration principle. On the other hand, suffixes constitute a larger set than prefixes. They come in two varieties, inflectional and derivational, both subcategories being quite large for a closed-set, and both exhibiting considerable allomorphy, as 85% of suffixes have allomorphs. The allomorphic changes apply to both stems and suffixes. More specifically, items sharing the same morphological (noun, verb or adjective, inflectional endings) and phonological features (same final character) exhibit similar allomorphic behavior.

- (4) a. εισφορά εισέρχομαι prefix: ΕΙΣ allomorph: –
(isfora) (iserxome)

‘contribution’	υπόλογος υφυπουργός (ipologos) (ifipurγos)	‘enter’ prefix: isallomorph: – prefix: ΥΠΟ allomorph: ΥΦ
‘accountable’	‘vice minister’ αντιμέτωπος ανθυγιεινός (antimetopos) (anθiyiinos)	prefix: ipo allomorph: if prefix: ANTI allomorph: ANΘ
‘opposing’	‘unhealthy’ μετατρέπω μεθεόρτια (metatrepo) (meθεortia)	prefix: anti allomorph: anth prefix: META allomorph: MEΘ
‘convert’	‘events after a feast’	prefix: meta allomorph: meth
b. ταξιτζής	ταξιτζήδες (taksitzis) (taksitziδes)	suffix: TZH(Σ) allomorph: TZHA
‘taxi driver’	‘taxi drivers’ παρκάρω παρκάρισα (parkaro) (parkarisa)	suffix: dzi(s) allomorph: dzidh suffix: AP(Ω) allomorph: API
‘I park’	‘I parked’ αβρότητα αβρότητες (avrotita) (avrotites)	suffix: ar(o) allomorph: ari suffix: OTHTA allomorph: OTHT
‘courtesy’	‘courtesies’	suffix: otita allomorph: otit

Melissaropoulou & Ralli (2009) deal with the general principles_which underlie the structural combination of a base with a particular suffix in Standard Modern Greek and some of its dialects. They argue that: a) suffixes select bases of a specific type, b) certain suffixes can be followed by other suffixes, while others are not susceptible to further suffixation, and c) the overall number of attested suffix combinations is generally smaller than the one theoretically possible.

The first systematic attempt to account for the combinatorial behavior of affixes was made within the framework of strata-oriented models (cf. Siegel 1974, Allen 1978, Selkirk 1982, Kiparsky 1982 and Mohanan 1986), according to which the different combinatorial properties of derivational affixes follow, to a great extent, from the position they hold into the different ‘lexical strata’ (‘levels’ in Kiparsky’s 1982 terms).

Therefore, in the light of evidence provided above, we argue in favor of the main thesis taken by Fabb (1988), Scalise (1994) and Melissaropoulou & Ralli (2009), according to which suffix-driven selectional restrictions are the ones that govern the formation of derivational structures.

3.3. Bound stems

Another case of interest in the morphological parsing of MG is a special type of words containing bound stems. As discussed in Petropoulou (2009) in this volume, this class of words comprises part of what we call neoclassical compounds in MG, because, like neoclassical compounds in English, they contain a bound element of Ancient Greek origin. Examples are *νηπι-αγωγ(ος)* (nipiayoyos) ‘preschool teacher’,

παθο-γον(ος) (pathogonos) ‘pathogenic’, *δακτυλο-γραφ(ος)* (daktilygrafos) ‘typist’, *σκηνο-θετη(ς)* (skinothetis) ‘director’, *τυρο-κομ(ος)* (tirokomos) ‘cheese producer’, *εντομο-κτον(ο)* (entomoktono) ‘insecticide’, *μετεωρο-λογ(ος)* (meteorologos) ‘meteorologist’, *καρδιο-παθ(η)ς* (kardiopathis) ‘cardiopath’, where the elements –*αγωγ(ος)* (–agogos), –*γον(ος)* (–gonos), –*γραφ(ος)* (–grafos), –*θετη(ς)* (–thetis), –*κομ(ος)* (–komos), –*κτον(ο)* (–ktono), –*λογ(ος)* (–logos) and –*παθ(η)ς* (–pathis) are bound morphemes, that is they cannot stand as free words.

Petropoulou (2009) has discussed the morphological status of these bound elements and the different opinions expressed which we present here in short. According to Giannouloupoulou (2000), following Anastasiadi–Simeonidi (1986), these elements are considered as ‘confixes’ (Martinet 1979), as they appear to acquire gradually more and more characteristics of suffixes. In these terms, confixes are secreted parts of words (Jespersen 1941, Warren 1990), which have been associated with a new specialized meaning. Examples of confixes cited by Giannouloupoulou (2000), are presented here with their extended meanings, such as –*λόγος* ((–logos) ‘scientist’ as above), –*λογία* ((–logia) ‘science’, as in *θεολογία* (theologia) ‘theology’), –*γράφος* ((–grafos) ‘writer/recorder’ as above), –*γραφία* ((–grafia) ‘science/study’, as in *ωκεανογραφία* (okeanografia) ‘oceanography’), –*κτόνος* ((–ktonos) ‘killer’, as above), –*κτονία* ((–ktonia) ‘killing’ as in *πατροκτονία* (patroktonia) ‘patricide’), –*ποιός* ((–pios) ‘maker’ as in *επιπλοποιός* (epiplopios) ‘carpenter/ (lit.) furniture maker’). For Giannouloupoulou, confixes constitute members of a closed set of items, which also includes initial elements such as *ευρω-* (evro–), *πολύ-* (poli–), *νέο-* (neo–), *παλαιο-* (paleo–), *τηλε-* (tile–) as well as the final bound element –*ισμός* ((–ismos) equivalent to the suffix –ism in English).

On the other hand, Ralli (2008a) supports that these elements are bound stems of a verbal origin and defies the argument favouring their suffixal character presenting a number of opposing arguments. She claims that: i) these elements can serve as bases to prefixation, e.g. *ιπο-λογος* (‘responsible for one’s actions’), *υπερ-μαχος* (‘supporter’), ii) they carry more concrete meaning in comparison to affixes which have a more functional role, often expressing agentive or instrumental meaning, iii) they carry valency information, i.e. information about the obligatory complements of the verbs they derive from, calling for theta–role saturation by the left–hand element in the constructions they appear, and iv) they participate in compound structures, which are recognizable both from the presence of the linking vowel –ο–, which constitutes a compound marker in Greek (Ralli 2008b), e.g. *πατρ-ο-κτονος* ((patroktonos) ‘patricide’ (agentive)) and from the recursivity they exhibit in their structures, e.g. *κοινωνι-ο-γλωσσ-ο-λόγος* ((kinoniologos) ‘socio–linguist’), which characterizes compounding.

The structures corresponding to the opposing views presented above for a word involving a bound element such as *βιολόγος* (viologos) ‘biologist’ are formulated as follows: a) *βιο-λογος*, where the element –*λογος* is a confix, and b) *βι-ο-λογ(ος)*, where the element –*λογ* is a bound stem. Although, there is seemingly no significant difference between the two structures, the implications they have for the

computational treatment of words containing these elements, are immense. This stems from the fact that as Ralli (2008a) has noticed, words containing bound elements, regularly serve as bases for the formation of derivatives, through suffixation, selecting suffixes from a closed set and giving rise to words such as *βιολογ-ια* (viologia) ‘biology’, *βιολογ-ικ(ος)* (viologikos) ‘biological’ and so forth. Confixation in this case, which renders the elements *-λογος* (–logos) and *-λογία* (–logia) as separate items belonging to the closed set of confixes, with no apparent morphological association between them, gives rise to the unrelated structures *βιο+–λογος* and *βιο+–λογία*, thus obscuring the obvious morphological relationship between the two first items. In these terms, the structure of the word *βιολόγος* (viologos) is not related to the structure of the word *βιολογία* (viologia), more than it is related, for example, to the structure of the word *βιογραφία* (viografia) ‘biography’ sharing with both of them only the same initial stem and a different confix. In computational terms, this would require the insertion of all possible confixes² (e.g. *-λογος* (–logos), *-λογία* (–logia), *-γράφος* (–grafos), *-γραφία* (–grafia), *-κτόνος* (–ktonos), *-κτονία* (–ktonia)) keeping them unrelated to each other. This would be quite inadequate as a morphological solution and not a very economical one for a computational analysis.

On the other hand, the ‘bound stem’ view gives rise to the structure *βι-ο-λογ(ος)*, which then, according to Ralli (2008a) serves as a base for the derivation of the word *βιολογία* (*βι-ο-λογ+ια*). In computational terms, this would require the insertion of all bound elements with verbal origin, along with the possible suffixes they may receive, namely the *-ια* (–ia), *-ικ-* (–ik–), *-ειο* (–io), *-ισσα* (–issa), *-ρια* (–ria), all of which are common suffixes in MG attaching to other bases apart from compounds with bound elements (e.g. *κατοικ-ια* (katikia) ‘residence’, *φιλ-ικ(ος)* (filikos) ‘friendly’, *Ασιάτισσα* (Asiatissa) ‘female Asian’ etc.). Apart from the obvious economy of the ‘bound stem’ solution, it serves for greater accuracy in the morphological analysis obtained, as it preserves the morphological relationships between words.

Therefore, supporting the ‘bound stem’ view, we compiled a corpus consisting of about 7000 words, each containing one of the 54 bound stems with verbal origin found in MG, such as *-λογ* (–log), *-γραφ* (–graf), *-κρατ* (–krat), *-δοτη* (–doti), *-δετη* (–deti), *-γον* (–gon), *-γεν* (–gen), *-μαθ* (–math), *-μαν* (–man) etc. along with their derivatives formed with the nominalising suffixes *-(e)ia*, *-(e)io*, *-issa*, *-ria* (e.g. *archeolog-ia* ‘archaeology’, *emodot-ria* ‘female blood donor’, *kosmogonia* ‘cosmogony’, *vivliodet-eio* ‘bookbinding site’) and verbs ending in (o) arising from conversion (e.g. *limokton(o)* ‘starve’).

4. The *Linguistica* Experiment

² The collection of confixes provided by Giannouloupoulou (2000) is not exhaustive, consisting only of a part of elements that could be classified as confixes, which may mean that potential confixes might have to satisfy a number of criteria in order to enter this class of items. This would leave out a significant number of elements, which would have to be treated in other terms.

Computational Analysis of Suffixes and Bound Stems of the Greek Language: A Crash Test with *Linguistica*

4.1. About *Linguistica*

Linguistica is a program designed to explore the unsupervised learning of natural language, with primary focus on morphology (word–structure). It runs under many operation systems, and is written in C++ within the Qt development framework. Its demands on memory depend on the size of the corpus being analyzed.

Unsupervised learning refers to the computational task of making inferences (and therefore acquiring knowledge) about the structure that lies behind some set of data, without any direct access to that structure. In the case of unsupervised learning of morphology, *Linguistica* explores the possibilities of morpheme–combinations for a set of words, based on no internal knowledge of the language from which the words are drawn.

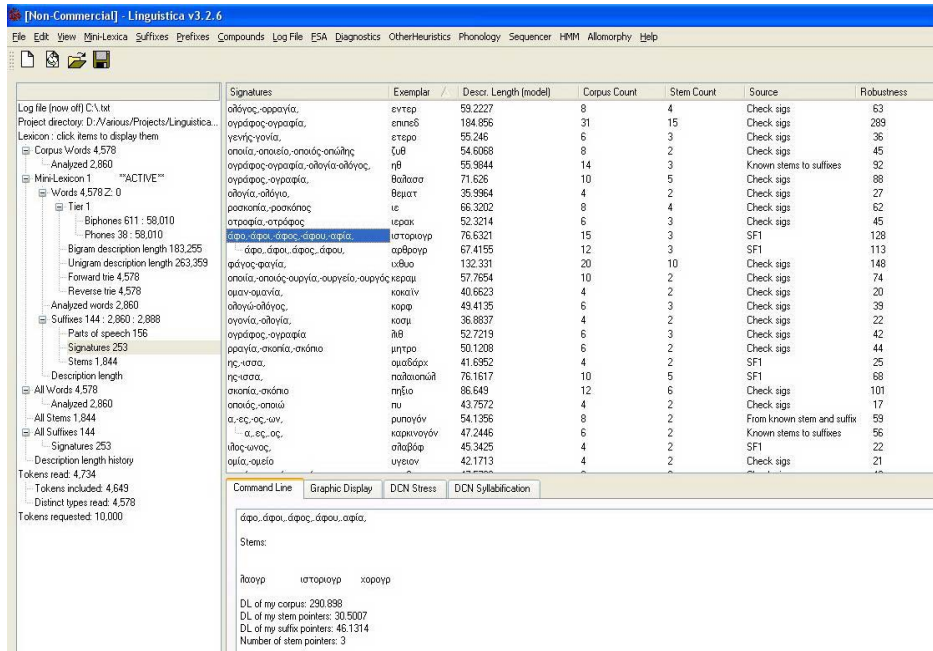


Figure 1: The interface of software *Linguistica*

Segmentation is the first task of this process; the program figures out where the morpheme boundaries are in the words, and then decides what the stems are, what the suffixes and so forth. Most of *Linguistica*'s functionality, at this point, goes into making these decisions. For our experiment, we used the 3.2.6 version (March 2009) for Windows XP.

4.2. Find Allomorphy with *Linguistica*

It is referred that *Linguistica* is capable of determining a limited amount of

allomorphy in stems. In many languages (including English), stem final material is deleted in front of certain suffixes. For example, stem-final *-e* is deleted in English before a number of suffixes: *love*, but *lov-ing* and not *love-ing*; *sane*, and *sanity*, not *sane-ity*. Goldsmith treats this as allomorphy, although it is not.

The strategy of *Linguistica* is to reanalyze material that had been previously included in a suffix as part of the stem, and provide the information that other suffixes must delete that material, when it appears before them. Goldsmith (2001) illustrates this with the following example: the words *love*, *loves*, *loved*, and *loving*, which had been analyzed as *lov* + signature *e.ed.es.ing*, will be reanalyzed with the stem *love* and the suffixes *NULL*, *ed*, *s*, and *ing*. The suffixes *-ed* and *-ing* will be informed that they are capable of deleting the preceding *e*, and this is indicated by placing an *e* in angle brackets before the prefix, thus: *<e>ing* and *<e>ed*. Thus the new signature for *love* is *NULL.<e>ed.<e>ing.s*, and this signature correctly deals both with stems that end in *-e* and those that do not.

Additionally it is pointed out that *Linguistica* treats *y*-final nouns and verbs in the same way: *academy/academies* are treated as if based on the stem *academy* and the suffixes *NULL* and *<y>ies*.

4.3. Our experiment corpora

As already put forward, our hypothesis is that *Linguistica* would appear to have major problems in analyzing a corpus of Greek words. In order to test this, three text corpora were created ad hoc; the first had 60,000 tokens (28,000 words) from a newspapers corpus, the second had 8,500 words with carefully selected lemmas and entries (words with same inflectional and derivational suffixes, groups of common prefixed words, etc) and the third was a science fiction novel with 200,000 words written by the first author of the present work. The results from the first corpus were quite disappointing (3% accuracy). The results from the third corpus were slightly better, but the accuracy was still quite low (6%). On the other hand, the results from the second corpus were more specific and clear, although the accuracy was also quite low. The system managed to detect several inflectional paradigms, few derivational suffixes and some bound stems. Additionally, only two allomorphy types were recognized, only one of which was correct, i.e. the *παιδι-παιδ*-type allomorphs!

4.4. Results

Checking our data with *Linguistica*, the top ten signatures are: (i.) *NULL.δεζ.δων*, (ii.) *ἀρειζ.ἀρετε.ἀρισα.ἀρονμε.ἀρω*, (iii.) *ἀρα.εζ.ηζ.ικός.ων*, (iv.) *άτων.ατάκι.ατάρα*, (v.) *ά.άκι.ου.ο.ων*, (vi.) *ά.άδεζ.άζ*, (vii.) *έζ.εδάκι*, (viii.) *ά.άζ.ατζή*, (ix.) *νεζ.νων* and (x.) *NULL.εις*. The first is composed of noun stems with a *δ*-allomorph (*μεζέ, κουβά, μαμά*), the second includes foreign stems, which form verbs with *-άρ(ω)* (*σκορ, σοκ, σκαν*) and the fifth is only combined with neutral nouns belonging to the sixth inflectional class (*βουνό, μωρό, νερό*).

The results are derived by the application of an advanced system with heuristics (see Goldsmith 2001). He Goldsmith points out that the overall sketch of the

morphology of English and other European languages comes out quite normal in its outlines. Nevertheless, the results from the English experiments, when studied closely, show that there are some parsing errors. The author of *Linguistica* tried quite successfully to fix these errors with additional heuristics and evaluate them using the MLD measure. However, the results from the Greek corpora do not require a closer study, since the errors form the rule rather than the exception. These errors may be organized in the following ways:

(a) The collapsing of two or more suffixes into one: for example, here we find the suffix *-ικός* (*-ikos*); in most corpora, the equally spurious suffix *-επτικός* (*-eftikos*) is found.

(b) The systematic inclusion of stem-final material into a set of (spurious) suffixes. In Greek, for example, the high frequency of stem-final *-τ* (*κύματ-α* (*kimata*)) can lead the system to the analysis of a set of suffixes as in the spurious signature *τος,τα.των* or *τακι.ταρα*.

(c) The inclusion of spurious signatures, largely derived from short stems and short suffixes, and the question related to the extent of the inclusion of signatures based on real, but overapplied, suffixes. For example, *-ς* (*-s*) is a real suffix of Greek, but not every word ending in *-ους* (*-us*) should be analyzed as contained that suffix.

(d) The failure to segment all words actually containing the same stem in a consistent fashion: for example, the stem *χορ* with the signature *ος.οι.ους* is not related to *χορ* with the signature *ενω.ενεις.ενει*.etc.

(e) Stems may be related in a language without being identical. The stem *αιμ* may be identified as appearing with the signature *α.ατα.ατο* and the stem *αι* may be identified with the signature *ματακι.ματαρα*, but these stems should be morphologically related.

(f) The system has never identified the linking vowel *-ο-* of the bound stems as a separate element. It was always attached either to the first component (*γλωσσο-*) or to the bound stem (*-ολόγος*) without any systematic treatment.

(g) *Linguistica* failed to treat correctly the allomorphy.

1. NULL.δες.δων	3. άρα.ές.ής.ικός.ών
αγά	αδερφ
βεζίρη	αυλ
γιαγιά	γραμμ
ζαρίφη	εποχ
ζελέ	μορφ
ζουρνατζή	φων
καναπέ	ψυχ
καυγατζή	4. άτων.ατάκι.ατάρα
καφέ	αιμ
κεσέ	αλμ
(...) 30 stems	αρμ
2. άρεις.άρετε.άρισα.άρουμε.άρω	βημ
κοπι	λημμ
παρκ	(...) 16 stems
σκαν	5. ά.άκι.ού.ό.ών
σκορ	βουν
σοκ	γλυκ
τρατ	μωρ
	νερ
	ποσ
	ποτ

Table 1: The top-five signatures of our second corpus

4.4.1. Prefixation

The analysis of prefixes in Greek should not pose a serious problem for *Linguistica*, since there are very few and with limited allomorphy. It managed to create signatures like *συν.αντι* *sin.anti* {εργατικός (*eryatikos*), ένζυμο (*enzimo*), εισφορά (*isfora*)}, *αντι.κατα* (*anti.kata*) {βάλλω (*valo*), θέτω (*theto*)}, *συν* (*sin*) {θετώ (*theto*, τρέχω (*trexo*), άγω(*ago*)}, which contain true prefixes. Nevertheless, as we mentioned before, signatures with two prefixes combined were also created, such as *συν.συνεπι* (*sin.sinepi*) {τηρω (*tiro*), τηρητής (*tiritis*), τηρούμαι (*tirume*)}, *συν.συνυπο* (*sin.sinipo*) {δηλώνω (*dilono*), δηλωτικός (*dilonotikos*)} and *αντι.συνυπο* (*anti.sinipo*) {γράφομαι (*grafome*), γεγραμμένος(*gegramenos*)}. Additionally, the system failed to relate prefixes with common characters like *α-* (*a-*) and *αν-* (*an-*), *κατα-* (*kata-*) and *κατ-* (*kat-*) or the most changeable prefix *συν-* (*sin*) {*συμ-* (*sim-*), *συγ-* (*sig-*), *σιλ-* (*sil-*), *συρ-* (*sir-*), *συσ-* (*sis-*)}, since the system does not incorporate any phonological rules, such as deletion and assimilation. Moreover, it was very common in spurious signatures to include some of the first characters of the stem in the prefixes (i.e. *συνδ-* (*sinδ-*), *συναρ-* (*sinar-*), *συνθηκ-* (*sinthik-*), *συναρμ-* (*sinarm-*)) or to mislabel part of stems as prefixes (*γλ-* (*gl-*), *λεν-* (*lef-*)). Finally, *Linguistica* could not detect any allomorphic behaviour of prefixes and of course it failed to relate them with other true forms of the same prefix, for example *κατα-* (*kata-*) and *καθ-* (*kath-*), *υπο-* (*ipo-*) and *υφ-* (*if-*).

Computational Analysis of Suffixes and Bound Stems of the Greek Language: A Crash Test with Linguistica

Signatures /	Mini	Exemplar	Desc. Length	Corpus Count	Sigs	Remarks
συν.συνε	δένω,	22.7458	4	2	PF1	13
συνε.συνεκ	φέρω,	17.3091	2	1	PF1	5
αντι.συνεκ	φωνώ,	22.9968	10	4	Singleton	53
αντικ.αντιχ	τυπέμαι,	23.5758	2	1	PF1	9
αντι	αισθητικά,	331.247	2580	1303	Singleton	5252
συνα.συναν	θροιά,	21.7685	2	1	PF1	7
συναρ.συνα	μοιλογύμαι,	21.4059	2	1	PF1	11
αντι.συνε	στραμμένος,	22.6489	4	2	PF1	26
NULL.αντι	συναδελφικά,	120.135	28	14	PF1	235
συν.συντ	αυτίζω,	31.2409	6	3	PF1	36
αντιφ	λογιστικός,	16.3615	1	1	PF1	0
συν.συνδι	αϊθάσσομαι,	56.3346	12	6	PF1	98
συνε	ιστά,	31.8486	3	3	PF1	8
αντικ	ιστά,	16.3615	1	1	PF1	0
αντιο.αντισ	υγκικός,	21.7685	2	1	PF1	8
αντι.αντιφ	λογιά,	17.3091	2	1	PF1	6
συνθέ	τιδα,	17.3615	1	1	PF1	0
συνα.συνυ	ποχωρώ,	21.7685	2	1	PF1	7
συν.συνεκ	τελεστής,	14.9465	2	1	PF1	9
συν.συντε	βήμιένος,	24.9682	4	2	PF1	28
συναρ.συνδι	αιτητής,	19.7685	2	1	PF1	8
αντι.συνα	μήτρια,	23.2863	4	2	PF1	25
συνύφ	ανση,	17.3615	1	1	PF1	0
συνδη.συνω	μάτσα,	24.5758	2	1	PF1	8
συν.συνδ.συντ	αυτίζω,	20.9909	3	1	PF1	14
συνα.συνε	γείρομαι,	26.1083	4	2	PF1	23
συνδ.	άνισα,	16.3615	1	1	PF1	0

Figure 2: Sample of prefix signatures of our corpus

4.4.2. Suffixation

The suffixal system of the Greek language is quite complex; as Melissaropoulou (2007a, 2007b) and Melissaropoulou & Ralli (2008) show, a stem can be followed by up to six suffixes (derivational and inflectional). *Linguistica* succeeded in creating some inflectional paradigms like the verbal present *ω.εις.ει.ουμε.ουτε.ουν* (γράφω (grafo) ‘write’, τρέχω (trexo) ‘run’) and *ο.ου.ων.α.[ακι]* (βουνό (vuno) ‘mountain’, μωρό (moro) ‘baby’, νερό (nero) ‘water’). Except for three other signatures, the rest of them (62) were spurious. There is an average number of signatures with combined suffixes (usually a derivational with an inflectional), such as *άρεις.αρει.αρω.αρουμε.αρισα, ατζη.ατζης.ατζηδες.ατζηδων* or *ευτικός.ευτικοί* (χορός (xoros) ‘dance’, δήμος (dimos) ‘municipality’). It was a very common mistake to create suffixes by including the last character of the stem; for example *γα.ζα* (ανοι (ani), τυλι (tili), διαλε (diale)) or *ινα.να* (γλυκα (glyka), πικρα (pikra), λευκα (lefka)).

Ω.ΕΙΣ.ΕΙ.ΟΥΜΕ.ΕΤΕ.ΟΥΝ
γραφ, τρεχ, δεν, βαζ, καν
Ε.ΙΝΟΣ.Ο.ΟΙ.ΟΣ.ΟΥ.ΟΥΣ.ΩΝ
ανθρωπ, κακτ, βαλτ
Ο.ΟΥ.Α.ΩΝ.ΑΚΙ
βουν,νερ, μωρ, κακ, ποτ
ΆΝΘΗΚΑ.ΆΝΘΗΚΕΣ.ΑΪΝΟΜΑΙ.ΑΪΝΟΥΜΕ.ΑΪΝΩ.ΑΝΘΕΪΣ.ΑΝΘΪ
λευκ, γλυκ, μωρ
ΓΑ.ΞΑ
ανοι, διαλε, κοιτα, τυλι
ΡΙΟΥ.ΡΙΩΝ.ΡΙ
καλαμα, ποτη, σαμα, σφν

Table 2: Signatures of inflectional and derivational suffixes

Goldsmith tried to fix this problem by advancing the heuristics and applying the feature “short-length for non-stems”; however, the treatment of one-character suffixes and prefixes is an important issue that causes many difficulties for a UML system. Finally, as claimed in our hypothesis, *Linguistica* failed to detect suffixal allomorphy, since the system did not relate the suffixes and usually failed to analyze them (45% failure). Therefore, it identified suffixes such as *αρω.αρισα* instead of *αρ~αρι* (*αρ<ι>*), *ατζής.ατζήδων.ατζήδες* instead of *τζη~τζηδ* (*τζη<δ>*) etc. As we can see, the accuracy of the system was 13% for derivational suffixes and 22% for inflectional suffixes³.

4.4.3. Stems

Linguistica presented a common behaviour in the analysis of nominal stems. First of all, only nominal allomorphs of the *παιδί*-type were detected. In the other cases, if there was a V-deletion allomorphy (i.e. *καρδιά~καρδι* (*karḗia~karḗi*) ‘heart’), the system detected only the V-deleted stem (*καρδι-*) considering the deleted vowel as a suffix. Moreover, if there was a C-insertion allomorphy (i.e. *κύμα~κυματ* (*kima~kimat*) ‘wave’), the system considered the final consonant of the allomorphs as the initial of the suffixes (*κύμα*). Additionally, there were a few signatures with spurious suffixes that contained the last two characters of the stem, such as *νας.να.νες.νων* (*σολη* (*solī*), *πυρη* (*pirī*), *αιω* (*eo*), *λιμε* (*lime*)) and *γα.ξα* (*ανοι* (*ani*), *τυλι* (*tili*), *διαλε* (*ḗiale*)). The system failed to relate any of the stems. Also the statistical analysis of both corpora reveals that only 4% of the allomorphs were detected by *Linguistica*. These results are similar to those of Kurimo *et al* (2007) for Finnish and Turkish; moreover, the hypothesis of *Linguistica*’s failure to deal with Greek allomorphy expressed by Karasimos (2008) was experimentally tested and

³ We consider as true signatures, the signatures that contain real suffixes. Of course, some signatures did not contain all the inflectional paradigms of a noun or a verb.

found to be valid.

Signatures	Mini	Exemplar	Desc.	Length	Corpus Count	Sigs	Remarks
ά.άδες.άς	αίλογ	26.731	3	1	1	Check sigs	8
γα.ξα	άνοι	43.5492	8	4	4	SF1	29
ικός.ών	αντρ	41.7779	8	4	4	SF1	32
ινος.ου.ών	αργύρ	40.5731	9	3	3	SF1	42
ου.ων	άρθρ	34.7184	6	3	3	SF1	22
ινος.ων.ου.ου	αρκούδ	36.5525	3	1	1	Known stems to suffixes	18
ίζω.ιστή.ιστής	βαπτ	48.4072	12	4	4	SF1	76
ικός	βασυή	63.8619	7	7	7	SF1	24
ικος	βεζίρ	57.0721	6	6	6	SF1	20
ματάκι.ματάρα	βη	43.5492	8	4	4	Check sigs	46
ά.άδες.άς.αδάκι	βορι	46.0905	12	3	3	Known stems to suffixes	63
ευτής.ευτικός	βουή	88.0206	22	11	11	Check sigs	163
ά.άς.ές.νός.ών	βραδι	42.6864	5	1	1	Loose fit	20
άδες.άς.ά.ά	γαϊτα	40.8237	6	2	2	Known stems to suffixes	41
ά.άδες.άς.ιού	γιαουρτ	37.0285	4	1	1	Loose fit	21
άρα.ές.ικός.ών	γραμμ	51.1773	16	4	4	SF1	78
ά.άς.ές.ών	γωνι	43.9105	12	3	3	From known stem and suffix	53
είς.ικός	εκδρομ	29.2365	4	2	2	SF1	20
ατζήδες.ατζής	ετοιμ	23.6773	2	1	1	Loose fit	5
ές.αδάκι	ζεή	48.2158	10	5	5	Known stems to suffixes	49
ά.άς.ατζήδες.ατζής	ζουρν	52.709	16	4	4	From known stem and suffix	96
άκι.ίζω.ιστή.ιστής	θερ	37.4848	4	1	1	SF1	9
ανθείς.ανθώ	θερμ	49.7152	10	5	5	Check sigs	58
αίνα.ανά	θέρμ	55.999	12	6	6	Check sigs	56
ά.άς.αδάκι	κουβ	28.0923	3	1	1	Loose fit	8
ά.άς.άτων.αδάκι	κυμ	35.8372	4	1	1	Loose fit	9
ά.άς.ατά.ατάκι	ζευκαυμ	42.1972	8	2	2	Loose fit	54

Figure 3: Signatures of nouns, verbs and adjectives

4.4.4. Compounds and Bound stems

As we already mentioned the inability to feed the system with any rules or structural information means that, despite our preferred morphological analysis of the words involving bound elements, the analysis obtained by the system would not necessarily be the desired one, which was indeed the case. Specifically, among the signatures produced by the analysis of our ‘bound–stem corpus’, we found the ‘real’ suffixes, such as the derivational *–ία* (e.g. *θεολογ–ία* (theologia) ‘theology’), *–είο* (e.g. *ανθοπωλ–είο* (anthopolio) ‘flower shop’), *–τη(ς)* (e.g. *αιμοδό–της* (emodotis) ‘blood donor’), *–ισσα* (e.g. *παλαιοπόλ–ισσα* (paleopolissa) ‘female antique seller’), the nominal inflectional (*ος*) (e.g. *βοτανολόγ(ος)* (votanologos) ‘votanologist’), (*ης*) (e.g. *πατριάρχ(ης)* (patriarchis) ‘patriarch’), and the verbal inflectional (*ω*) (e.g. *ηχογραφ(ώ)* (ixografo) ‘sound record’). However, we also found sequences like *–ολόγος* (*–ologos*), *–ολογία* (*–ologia*), *–ογράφος* (*–ografos*), *–ογραφία* (*–ografia*), *–ομανής* (*–omanis*), *–ομανία* (*–omania*), *–οποιία* (*–opiia*), *–οποιείο* (*–opiio*), *–οτρόφος* (*–otrofos*), *–οτροφία* (*–otrofia*), *–όφιλος* (*–ofilos*), *–ορραγία* (*–orragia*), *–ογονία* (*–ogonia*), *–οστάτης* (*–ostatis*), *–οφαγία* (*–ofagia*), *–οκτονία* (*–oktonia*), which are basically like confixes with the linking element attached to them. At the same time, and for no obvious reason, among the signatures, we found sequences like *–φάγος* (*–fagos*), *–φαγία* (*–fagia*), *–παθής* (*–patis*), *–σκοπία* (*–skopia*), *–σκόπιο* (*–skopio*), *–ούχος* (*–uxos*), *–γενής* (*–genis*), *–γονία* (*–gonia*), *–μαθής* (*–matis*), *–άρχης* (*–axis*), *–φόρος* (*–foros*), *–πρεπής* (*–prepis*), *–τέχνης* (*–texnis*), which are also confix–like but without the element *–ο–* attached to them. Results like these, imply that the system

did not manage to recognize neither the linking element as a separate entity, nor the derivational or inflectional suffixes attached to the final bound elements.

Reasonably enough, the recognition of a great number of confix-like sequences with the linking element attached, as those mentioned above, gave rise to a great number of ‘correct’ stems⁴ of MG like *miθ-* (‘myth’), *okean-* (‘ocean’), *selin-* (‘moon’), *musik-* (‘music’), *xart-* (‘paper’), *sidir-* (‘iron’) or stem allomorphs like *δramat-* (‘drama’), *xromat-* (‘colour’), *θavmat-* (‘miracle’), *stromat-* (‘mattress’), *nimat-* (‘thread’) and so on, appearing as right hand elements in the words provided. However, also as stems were recognized sequences that are like compound stems, such as *kriptograf-*, *sismology-*, *karkinoγon-*, *vivlioklop-*, *texnolog-*, *plutokrat-*, due to the recognition of true derivational and inflectional suffixes that we saw above.

As a conclusion, we should note that the system did not manage to recognize any of the bound stems such as *-logy*, *-graf*, *-kton*, *-math*, *-krat* and so on, neither the linking element *-o-*, as proposed by the preferred morphological analysis for the words involving bound elements in MG. As we mentioned above, this fact was basically due to the lack of any morphological input to the system, which could lead the morphological analysis towards a particular direction.

Linguistica could not analyze any compounds. Its strategy and architecture is to extract suffixes and prefixes even for languages with rich morphology. English corpora that were tested in this system contained very few one-word compounds and a significant group of neoclassical compounds; the authors do not show that this system treated them correctly. Unfortunately the three Greek test corpora cannot serve as the basis for any serious conclusions for Greek compounds, since the results were totally haphazard. As a rule, the system was unable to recognize any of the compound’s components and failed to analyze many of them.

Stem	Phonological content	Length ptr to me	Corpus count	Suffix sig
κερ	15.6284	12.4888	1	εοκόςκος
κπρ	15.6284	10.9039	3	οποιός, οποιός, οσάτής.
κρυ	15.6284	11.4888	2	ολόγησα...οσάτής.
λεξ	15.6284	12.4888	1	ιλόγιο.
ληθ	15.6284	11.4888	2	ογράφος...ογραφία
ληπ	15.6284	11.4888	2	ομην.ομανία.
μην	15.6284	11.4888	2	ολόγιο...ορραγία.
μμ	15.6284	11.4888	2	ογράφος.ογραφία.
μν	15.6284	10.4888	4	όρρη.ογραφία...ομανής.ομανία.
μυθ	15.6284	9.90388	6	ογράφος.ολόγος...ομανής.ομανία...οποιός
μυρ	15.6284	11.4888	2	ολόγιο.ολόγος.
νεκ	15.6284	11.4888	2	ροσκόπιο...ροσκόπος
νεφ	15.6284	12.4888	1	ολόγιο.
νομ	15.6284	10.4888	4	όρρη.ογραφία...ολόγιο...ολόγος
νοσ	15.6284	10.1669	5	ογραφία...ολόγιο.ολόγος...ομανής.ομανία.
ξεν	15.6284	11.4888	2	ομανής.ομανία.
ξυλ	15.6284	10.1669	5	ογράφος...οσάτής...ουργία...ουργεία...ουργός
ογκ	15.6284	12.4888	1	ολόγος.
οκν	15.6284	12.4888	1	ολόγος.
οθ	15.6284	12.4888	1	οσής.
οντ	15.6284	11.4888	2	ογονία...ολόγος.
ορε	15.6284	11.4888	2	ιβάτης...ογονία.
οσμ	15.6284	12.4888	1	ολόγιο.
ουλ	15.6284	12.4888	1	ορραγία.
ουρ	15.6284	12.4888	1	ολόγος.
οσμ	15.6284	11.4888	2	φάγος.φαγία.

Figure 4: Signatures of bound stems

⁴ i.e. without their inflectional ending as they normally appear in compounds.

On the other hand, the results from Greek bound stems (neoclassical stems) were quite enlightening. We tested a corpus using more than 7,000 examples from a corpus by Petropoulou. *Linguistica* created many signatures with true bound stems; for example *λόγος* (*ακριβο* (akrivo), *επιγραφο* (epigrafo), *γλωσσο* (γλοso)), *αρχης* (*arhis*) *γεν* (gen), *γυμνασι* (gimnasi)), *μαθής.μαθεία* (mathis.mathia), *ελληνο* (elino), *αγγλο* (aglo)), *μανής.μανία* (manis.mania) *διψο* (dipso), *δοξο* (dokso), *ξενο* (kseno)). Studying these results in greater detail, it becomes obvious that the system scored better with the bound stems. Nevertheless, the linking vowel *-o-* was unpredictably attached either to the first component (*γλωσσο* (γλοso), *ελληνο* (elino)) or to the second component (*-ολόγος* (-ologos), *-ογραφία* (-ografia)); it was never analyzed as a separate element of these words. Additionally almost all inflectional suffixes were segmented as part of the bound stems

5. Conclusions

Computational Morphology is a rapidly growing area of linguistics. Unsupervised Morphology Learning Theory is a recent approach to morphological analysis problems, and seems to work well for languages with poor inflectional morphology, although any attempt to use this theory in morphologically rich languages, such as Finnish and Turkish, could be characterized at least as mediocre (Kurimo et al. 2006, 2008). We claim that a system without: a.) prior human-designed analysis of the grammatical morphemes of a language, b.) some identifying stems and affixes and c.) pre-imported morphological and phonological rules for correct parsing, is bound to fail. A system which builds lexica based on a common sequence of phonemes without proper rules is unable to treat successfully the complex combinations/behaviour of derivational suffixes and bound stems. As already shown, the phenomenon of allomorphy in Greek is very extensive. Allomorphy participates with the same frequency in every word formation process. A natural question to ask is whether a UML model is able to analyze processes and successfully treat suffixes and bound stems. We have presented a considerable amount of data with allomorphs and shown the complexity of the allomorphic changes, the combinability of derivational affixes and the normality of bound stems. Since the insertion of processing rules for allomorphy is not allowed in a UML model, the goal of correct parsing will never be attained. From a more theoretical point of view, our work has nothing to do with the current question: does a young speaker learn a language and segment the morphemes the way that a UML does? Thus, we would like to point out that only supervised morphology learning models with rules and imported human knowledge can serve as the basis for the computational treatment of the morphological phenomena of derivation and compounding in Modern Greek.

Bibliography

Allen, J., Hunnicutt, M. S. & D. Klatt (1987). *From Text to Speech: the MITalk System*.

- Cambridge: Cambridge University Press.
- Anastasiadi-Symeonidi, A. (1986). *Neology in Common Modern Greek* [In Greek]. Thessaloniki.
- Andron, D. (1962). *Analyse morphologique du substantif russe*, Tech.rep., Centre d' Etudes pour la Traduction Automatique, Universite de Grenoble 1.
- Bernard-Georges, A., Laurent, G., & D. Levenbach (1962). *Analuse morphologique du verbe allemande*. Tech. rep., Centre d' Etudes pour la Traduction Automatique, Universite de Grenoble1.
- Booij, G. (1997). Allomorphy and the Autonomy of Morphology, *Folia Linguistica XXXI/1-2*: 25–56.
- Booij, G. (2002). *The morphology of Dutch*. Oxford/ New York: Oxford University Press.
- Boussard, A., & M. Berthaud (1965). *Presentation de la synthese morphologique du français*. Tech. rep., Centre d' Etudes pour la Traduction Automatique, Universite de Grenoble1.
- Brand, I., Klimonow, G. & S. Nündel (1969). *Lexiko-morphologische Analyse*. In: Nündel, S., Klimonow, G., Starke, I., and Brand, I. (Eds.), *Automatische Sprachübersetzung: Russisch-deutsch*, 22–64. Akademie-Verlag, Berlin.
- Carstairs, A., 1989. *Allomorphy in inflection*. London: Croom Helm.
- Creutz, M. & K. Lagus (2002). Unsupervised discovery of morphemes. *Proceedings of the ACL-02 Workshop on Morphological and Phonological Learning*, 21–30.
- Drachman, G. (2006). A note on 'shared' allomorphs. *Journal of Greek Linguistics 2006*: 5–37.
- Dzeroski, S. & E. Tomaz (1997). Induction of Slovene nominal paradigms. In Nada Lavrac and Saso Dzeroski, (eds), *Inductive Logic Programming, 7th International Workshop, ILP-97*: 17–20. Lecture Notes in Computer Science, Vol. 1297. Prague, Czech Republic. Springer: Berlin.
- Flenner, G. (1995). Quantitative Morphsegmentierung im Spanischen auf phonologischer Basis. *Sprache und Datenverarbeitung*, 19(2): 63–79.
- Galani, A. (2003). Allomorphy: Theme Vowels in Modern Greek. *Camling Proceedings 1*: 100–107.
- Giannouloupoulou, G. (2000). *Morphosemantic comparison of Affixes and Confixes in Modern Greek and Italian* [In Greek]. Thessaloniki.
- Goldsmith, J. (2001). Unsupervised Learning of the Morphology of a Natural Language, *Computational Linguistics 27*, vol 2: 153–196.
- Hafer, M. & S. Weiss (1974). Word segmentation by letter successor varieties. *Information Storage and Retrieval 10*: 371–385.
- Harris, Z. (1955). From phoneme to morpheme. *Language*, 31: 190–222. (Reprinted in 1970).
- Harris, Z. (1967). Morpheme boundaries within words: Report on a computer test. *Transformations and Discourse Analysis Papers 73*. Department of Linguistics, University of Pennsylvania.
- Hutchins, W. J. (2001). Machine translation over 50 years. *Histoire, Epistemologie Language 22 (1)*: 7–31.
- Janssen, A. (1992). Segmentierung französischer Wortformen in Morphe ohne Verwendung eines Lexikons. In Ursula Klenk (ed.), *Computatio Linguae*, 74–95. Steiner Verlag, Stuttgart.
- Kaplan R. M. & M. Kay (1981). *Phonological rules and Finite State Transducers*. Winter Meeting of ACL/LSA 1981.
- Karttunen L. (1983). KIMMO: A general morphological processor. *Texas Linguistics Forum 22*: 163–186.

Computational Analysis of Suffixes and Bound Stems of the Greek Language: A Crash Test with Linguistica

- Klenk, U. (1992). Verfahren morphologischer Segmentierung und die Wortstruktur im Spanischen. In Ursula Klenk (ed.), *Computatio Linguae*. Steiner Verlag, Stuttgart.
- Koskenniemi, K. (1983). Two-Level Morphology: A General Computational Model for Word-form Recognition and Production, *Proceedings COLING '84*, 178–181.
- Kurimo, M., Creutz, M., Varjokallio, M., Arisoy, E. & M. Saraclar (2006). Unsupervised segmentation of words into morphemes – Challenge 2005: An Introduction and Evaluation Report. *Journal of Proceedings ICSLP 2006*.
- Kurimo M., Mathias C. & M. Varjokallio (2008). Unsupervised Morpheme Analysis Evaluation by a Comparison to a Linguistic Gold Standard Morpho Challenge 2007 In A. Nardi & C. Peters (eds.) *Working Notes of the CLEF 2007 Workshop*.
- Lieber, R. (1982). Allomorphy. *Linguistics Analysis 10(1)*: 27–52.
- Mascaró, J. (1996). External allomorphy as emergence of the unmarked. In J. Durand & B. Laks (eds.) *Current Trends in Phonology: Models and Methods*: 473–483. Salford: European Studies Research Institute.
- Mascaró, J. (2007). External allomorphy and lexical representation. *Linguistic Inquiry 2007*.
- Martinet, A. (1979). *Grammaire Fonctionnelle du Français*. Paris: Didier.
- Matthews, P. (1965). A procedure for morphological encoding. *Mechanical Translation 9*: 15–21.
- McIlroy, M. D. (1982). Development of a spelling list. *IEEE Transactions of Communications 30 (1)*: 91–99.
- Melissaropoulou D. & A. Ralli (2008, submitted in Morphology). Structural combinatorial properties of Greek derivational suffixes. Paper presented at the 13th International Morphology Meeting (Vienna, February 3-6 2008), Workshop on Affix Ordering
- Melissaropoulou, D. (2007b, in print). Remarks on the combinability of derivational suffixes in Greek and its dialectal variation. In *Proceedings of the 3rd International Conference on Modern Greek Dialects and Linguistic Theory* (Nikosia, 14-16/06/2007).
- Porter, M. (1980). An algorithm for suffix stripping. *Program 14 (3)*: 130–137.
- Petropoulou, E. (2009). On the parallel between Neoclassical compounds in English and Modern Greek. In A. Ralli (ed.) *Patras Working Papers in Linguistics, Vol.1. Special Issue: Morphology*. Centre of Modern Greek Dialects. Department of Philology. University of Patras.
- Ralli, A. (1994). Feature Representations and Feature-Passing operations in Greek Nominal Inflection. *Proceedings of the 8th Symposium on English and Greek Linguistics*: 19–46. Thessaloniki: English Dept. Aristotle University of Thessaloniki.
- Ralli, A. (2000). A feature-based analysis of Greek nominal inflection, *Glossologia 11–12*: 201–228.
- Ralli, A. (2006). On the role of Allomorphy in inflectional Morphology: Evidence from Dialectal variation. *Advances of Language Studies 1*: 1–32.
- Ralli, A. (2008a). Greek Deverbal Compounds with Bound Stems. *Journal of Southern Linguistics 29 (1/2)*: 150-173.
- Ralli, A. (2008b). Compound Markers and Parametric Variation. *Sprachtypologie und Universalienforschung* (STUF) 61(1): 19-38.
- Ralli, A. (2009). Hellenic Compounding. In R. Lieber & P. Stekauer (eds.) *The Oxford Handbook of Compounds*, 453-464. Oxford: Oxford University Press.
- Ralli A. & E. Galiotou (1987). A Morphological Processor for Modern Greek. In *Proceedings of ACL (Association for Computational Linguistics)*: 26–31. Copenhagen.
- Ralli, A. & E. Galiotou (2001). A Prototype for a Computational Analysis of Modern Greek Compounds. *Asymmetry Conference '2001*. Montréal: UQAM.

- Ralli, A. & E. Galiotou (2003). Processing Greek Compounds. *Proceedings of the ESF/SCH Exploratory Workshop "Constructing Bilingual Computerised Dictionaries with Special Emphasis on Lesser Used Languages"*. Aristotle University, Thessaloniki.
- Ralli, A. & E. Galiotou (2005). Greek Compounds: A Challenging Case for the Parsing Techniques of PC-KIMMO v.2. *International Journal of Computational Intelligence*, vol. 1, no. 2: 128–138.
- Rissanen, J. (1989). *Stochastic Complexity in Statistical Inquiry*. Singapore: World Scientific Publishing Co.
- Roark B. & R. Sproat (2007). *Computational Approaches to Morphology and Syntax*. Oxford: Oxford University Press.
- Schone, P. & D. Jurafsky (2001). Knowledge-free induction of morphology using latent semantic analysis. *Proceedings of the 4th Conference on Computational Natural Language Learning (CoNLL)*: 67–72.
- Schveiger, P. & J. Mathe (1965). Analyse d' information de la declinaison du substantive en hongrois (du point de vue de la transduction automatique). *Cahiers de Linguistique Theorique et Appliquee* 2: 263–265.
- Thornton, A. (1997). Stem allomorphs, suffix allomorphs, interfixes or different suffixes? On Italian derivatives with antesuffixal glides. In G. Booij, S. Scalise & A. Ralli (eds.) *Proceedings of 1st Mediterranean Meeting of Morphology*, 78–91.
- Vauquois, B. (1965). *Presenetaion d' un programme d' analyse morphologique russe*. Tech. rep., Centre d' Etudes pour la Tranduction Automatique, Universite de Grenoble1.
- Warren, Beatrice (1990). The importance of combining forms. In Wolfgang Dressler et al. (eds) *Contemporary Morphology*. Berlin: Mouton de Gruyter.
- Woyna, A. (1962). *Morphological analysis of Polish verbs in terms of machine translation*. Tech. rep., Machine Translation Research Project, Georgetown University.
- Yarowsky, D. & R. Wicentowski (2001). Minimally supervised morphological analysis by multimodal alignment. *Proceeding of the 39th Annual Meeting of the Association for Computational Linguistics (ACL)*: 207–216.
- Καρασίμος, Α. (2001). *Η αλλομορφία στην κλίση και τη σύνθεση της Ελληνικής Γλώσσας*, Πτυχιακή εργασία, Πανεπιστήμιο Πατρών.
- Μαρκόπουλος, Γ. (1998). *Η επεξεργασία του ονόματος στα ελληνικά*, Διδακτορική διατριβή, Πανεπιστήμιο Αθηνών.
- Μελισσαροπούλου, Δ. (2007α). *Μορφολογική περιγραφή και ανάλυση του μικρασιάτικου ιδιώματος της περιοχής Κυδωνίων και Μοσχονησίων: η παραγωγή λέξεων*, Διδακτορική διατριβή, Πανεπιστήμιο Πατρών.
- Ράλλη, Α. (2005). *Μορφολογία*. Αθήνα: Πατάκη.
- Σγάρμπας, Κ. (1997). *Τεχνικές αυτόματης μορφοσυντακτικής ανάλυσης της Νέας Ελληνικής Γλώσσας*, Διδακτορική διατριβή, Πανεπιστήμιο Πατρών.

ATHANASIOS KARASIMOS
UNIVERSITY OF PATRAS
akarasimos@upatras.gr

EVANTHIA PETROPOULOU
UNIVERSITY OF PATRAS
evapetro@cc.uoi.gr