

**Enhancing a Greek Language Stemmer**

-

**Efficiency and Accuracy Improvements**

Spyridon Saroukos

University of Tampere  
Department of Computer Sciences  
Computer Science M.Sc. thesis  
Supervisor: Eleni Berki  
July 2008

University of Tampere

Department of Computer Sciences

Spyridon Saroukos: Enhancing a Greek Language Stemmer - Efficiency and Accuracy Improvements

M.Sc. thesis, 48 pages, 20 Appendix pages

July 2008

---

Stemming algorithms are used in the field of Information Retrieval in order to improve precision and recall. Although for Greek there are three stemmers published, only one of them is freely available. In this thesis, we use stemmer performance metrics for evaluating the existing algorithm and we improved its accuracy and completeness. These improvements were achieved by providing an alternative implementation in PHP which offers more syntactical rules and exceptions. Finally, the two algorithms are tested and their statistics metrics are compared.

Key words and terms: stemming algorithm, Greek, algorithmic efficiency, PHP

## Table of Contents

1. Introduction.....	1
1.1 Stemming.....	1
1.2 Definitions of Key Terms and Concepts.....	2
1.3 Inflectional versus Derivational Variants.....	3
1.4 Stemming Techniques – Advantages and Disadvantages.....	4
1.5 Greek Stemmers.....	4
1.6 Problems and Issues with the Latest Algorithm for Greek.....	6
1.7 Aims of this Thesis.....	7
1.8 Research Question.....	8
1.9 Methodology.....	8
1.10 Overview of Thesis' Contents.....	9
2. The Greek Language.....	10
2.1 The History of the Greek Language.....	10
2.2 Stemming in Greek.....	12
3. Stemmer Performance Metrics.....	14
3.1 Frakes' Metrics .....	14
3.2 Error Metrics.....	15
4. Algorithmic Design.....	17
4.1 The Existing Algorithm .....	17
4.2 Rejected Designs.....	24
4.2.1 Lovins' Design.....	24
4.2.2 Context-Free Grammar Design.....	25
4.2.3 Dictionary Based Design.....	25
4.2.4 Krovetz's Experimental Algorithm.....	26
4.3 Evaluation of Ntais' Algorithm.....	26
5. Our Improvements.....	28
5.1 Introduction of Stop-Word Elimination .....	28
5.2 Addition of More Grammatical Rules.....	28
5.3 Introduction of Lower Case Letters.....	29
5.4 The Revised Algorithm .....	29
6. Final Evaluation .....	42
7. Conclusion .....	44
References.....	46
Appendices.....	I
Appendix A: Verb Conjugation Classes In Greek.....	I
Verbs of 1st Conjugation Classes.....	I
Verbs of 2nd Conjugation Class.....	VII
Appendix B: Evaluation of Modified Algorithm Testing Output.....	XIII
Appendix B: On-line Stemmer .....	XVI
Appendix C: Stop Word List.....	XVII

## Index of Tables

Table 1: Stemming Examples.....	1
Table 2: Definitions of key terms in the stemming process .....	2
Table 3: Stemming Algorithms – Summarized Information.....	6
Table 4: The Greek Alphabet. Letters and their equivalent sound in English.....	11
Table 5: The noun “cat” in Greek.....	12
Table 6: Stemming Errors.....	15
Table 7: Ntais' Algorithm.....	17
Table 8: Our revised algorithm (additions and modifications highlighted).....	30
Table 9: Statistics: Comparison of the original and revised algorithm.....	42

## 1. Introduction

Search engines play a critical role in people's life nowadays. Google reported back in 2006 that it receives more than 100 million queries from US based hosts daily [Witten *et al.*, 2007]. For end users, they are perhaps the only way to find the information they seek and to navigate to the appropriate web pages. In addition, the users of search engines tend to navigate through the top results that the search engines return to them. It is quite common that end users know which piece and what type of information they seek but they are quite often unable to construct a proper query that will fully describe their request. As a result, malformed and badly structured queries tend to return fewer and more irrelevant results than well/structured queries. For example, someone seeking information for “World War II” may form a query as “World Wars”. Since “Wars” is not a part of “War”, documents containing the reference “War” will not be returned and, thus, the number of relevant results will be reduced. One of the solutions proposed to increase the ratio of relevant documents to total documents retrieved, also known as *recall*, is *stemming*.

### 1.1 Stemming

Stemming is the process of reducing words to their stem, base or root form, [Lovins, 1968] as shown in Table 1.

**Table 1: Stemming Examples**

<b>in English</b>	
<b>Original word</b>	<b>Stem</b>
Seas	sea
Wars	war
determination	determin-
Developed	develop
<b>in Greek</b>	
Αναπνέω	αναπνέ-
Ελληνικός	ελληνικ-

This is a simple but effective operation used in the fields of information extraction and natural languages [Carlberger *et al.*, 2001]. Stemming can be utilized when storing

information about a web page in a search engine's database or for *query expansion* [Xu & Croft, 1998]. In that case, the user's query is evaluated and reformulated. The search engine may reduce the word to its stem and thus return more results to the user. An evaluation research in 1981 showed that stemming improved search precision [Brants, 2003].

## 1.2 Definitions of Key Terms and Concepts

As stemming is also a linguistic process, any discussion of it, even in the field of IR, assumes the knowledge of some basic linguistic terms. The ones used in this thesis are described in Table 2.

**Table 2: Definitions of key terms in the stemming process**

Stem: a base part of a word that may have or may not have semantic meaning and no <i>affixes</i> (see below)
Affix: a small linguistic unit with semantic meaning that is attached to the beginning or the end of a stem to form a word
Prefix: an affix that is added to the beginning of a word or stem
Suffix: an affix that is added to the end of a word or stem
Stop Word: a term that appears so frequently in documents that it does not help searches [van Rijsbergen, 1979]. For example: <i>I, a, to, any, where, you</i>
Inflection: the modification or marking of a word to reflect semantic and grammatical information like gender, tense, number, case or person.  For example: to help - > <i>helps, helping, helped</i> (reflects gender and tense)
Derivation : the modification of a word that transforms it from one syntactical category (verb, noun, adverb, adjective) to another.  For example: to hope - > <i>hopefully, hopeless, hopes</i> (transformation to adverb, adjective and noun/verb)

Conflation class: a group of words that share the same semantic meaning [Paice, 1996].

For example:

*group , grouping, teams*

Compounding: the creation of new words by combination of two or more different words into a single form. For example:

*solar-powered, breastfeeding, bitter-sweet, antidisestablishmentarianism*

Morphological Variants: two or more words that are related due to inflection, derivation of compounding.

For example:

*dark, darkening, darks, darkroom*

### 1.3 Inflectional versus Derivational Variants

Inflectional and derivational are the two categories of morphological variants that stemmers mainly deal with. Any given word can have inflectional and derivational morphological variants. Inflectional morphological variants share the same basic meaning and belong to the same part of speech. For instance, the changes may affect the word's case, number, tense and gender. In contrast derivational morphological variants can belong to a different part of speech and thus, may mean something completely different from their stem since an adjective can be derived from a noun, or a noun from a verb among others. Many of the algorithms so far developed do not treat derivational suffixes or handle them partially. According to Paice [1994], affixes may contain important information about the meaning of the word and so it is advisable not to discard it during the stemming process. Stemming “antidote” to “dote “ creates a word that belongs to a different conflation class since the two words deal with different concepts.

Paice refers to the English language. However, the same phenomenon is noticed also in Greek. Affixes that are added usually alter the meaning of the word completely compared to the initial meaning of the stem. Some examples of this type of derivation are “μόρφωση” (education) to “παραμόρφωση” (disfigurement, deformation) and “δοκιμάζω” (I try) to “αποδοκιμάζω” (I boo, I abjure). A stemmer that includes derivational rules can be helpful for language research, but may be inappropriate as a query expansion tool, since it can supply the search engine with variants of the original

words that have a very different meaning. This fact can increase the number of matched documents but the semantic accuracy to the original word will be low.

#### **1.4 Stemming Techniques – Advantages and Disadvantages**

In addition to whether or not a stemmer treats prefixes, they are also categorized as (i) dictionary based (ii) based on algorithms or (iii) a hybrid version of both [Ntais, 2006]. Dictionary based stemmers use ready made dictionaries and match a word with its stem from a list. The main drawback of such stemmers is that dictionary maintenance is required and that these stemmers can not scale to handle unlimited words.

In contrast, algorithm based stemmers have been the focus of research with the algorithms of Lovins and Porter being the most representative. The former [Lovins, 1968] precedes the latter by 12 years and it was the first stemming algorithm ever published. It uses an extensive list of 294 endings, 35 transformation rules and 29 conditions. The algorithm is executed in two basic steps. In the first step the longest ending is matched and removed and in the second step the algorithm checks whether one of the 35 transformation rules should be applied.

Although Lovins' algorithm was the first stemming algorithm, Porter's is considered one of the most influential [Krovetz, 1993; Ntais, 2006; Frakes, 2003]. It was initially written for the English language and later ported to other European languages like Italian, Spanish, French and Portuguese. In contrast with Lovins' algorithm, it iteratively applies a set of rules and removes suffixes until no rules apply. The execution is completed into five distinct steps [Porter, 1980] and is considered very aggressive [Krovetz, 1993; Xu and Croft 1998]. Porter's work became influential and many implementations were written and made available by others. Unfortunately, these implementations contained errors. In order to deal with this problem, Porter released a framework with which stemming algorithms can be implemented using a string handling programming language called Snowball.

#### **1.5 Greek Stemmers**

Three stemmers have been developed for Greek. The first two are the TZK algorithm by Kalamboukis and Nikolaidis in 1995 and the Automated Morphological Processor (AMP) by Tambouratzis and Carayanis in 2001. These two stemmers have an acceptable accuracy of 90 to 95% but for their development some constraints were implied. The



AMP algorithm assumes that each word just consists of a stem part and an ending part and thus excludes all compound words. On the other hand, the TZK algorithm can only manipulate 65 suffixes although there are more than 166 suffixes in the Greek language. Furthermore, to our knowledge, neither of these algorithms has a freely available implementation.

The latest stemming algorithm developed for the Greek language is by George Ntais, in 2006. This algorithm follows the structure of the Porter algorithm and has a free implementation available on the web [Ntais, 2008]. The author has provided a web interface where users can make simple queries by posting a single word in Greek and have the word's stem returned. The interface is simple and uses Javascript for the implementation of the algorithm. According to the author, the algorithm can handle 158 suffixes of the Greek language, clearly outperforming the TZK and AMP algorithms. Nevertheless, in order to avoid complexity and due to constraints imposed during its development, it can only work with words in capital letters. In the Greek language, lower case words have accent marks that can totally change the meaning of a word, like the adjectives “αβαθής” and “ἄβαθης”. Both of the words mean “shallow”, with the former being the masculine nominative case and the later the feminine genitive. In addition, the stemmer is able to handle suffixes but not prefixes. The essential information about all stemming algorithms that affected this thesis work and were presented in the previous two sections is listed in Table 3.

**Table 3: Stemming Algorithms – Summarized Information**

Name	Language	Year	Web Availability	Derivation Dealing	Number of Execution Steps	Weaknesses
Lovins'	English	1968	Yes	Yes	2	Aggressive with short stems and words
Porter's	English	1980	Yes	No	5	Quite aggressive and produces overstemming [Carlberger <i>et al.</i> , 2001]
TZK	Greek	1995	No	Yes	2	Does not handle all suffixes
AMP	Greek	2001	No	No	4	Unable to handle compound words
Ntais'	Greek	2006	Yes	No	29	Relatively new and untested; can only handle capital letters; produces understemming errors

### 1.6 Problems and Issues with the Latest Algorithm for Greek

According to its author [Ntais, 2006], the latest stemming algorithm for Greek suffers from a few limitations and constraints that had to be imposed during its development. One of them is the incapability to handle lower case letters. The algorithm is only able to handle words in upper case letters and will not stem any word that contains even one letter in lower case. Since words in upper case do not have tone marks in Greek, the author is solving the problem of the moving mark phenomenon that can be observed during the conjugation of verbs, nouns and adjectives since no tone mark has to be presented in the stem returned. Because of that, the algorithm has a new limitation since most of the words encountered in Greek texts are in lower case.

Another issue with the existing algorithm is incapability to process some important

suffixes. The rationale was that the inclusion of these suffixes would introduce more errors if the appropriate exception list was also introduced. The creation of an extensive exception list was not feasible at that moment, so the algorithm is not treating suffixes like “-ατε” ,” -αστε” and “-τε”. These suffixes correspond to words from many syntactical categories like adverbs, nouns and verbs. Verbs with these endings correspond to past tenses. Past tenses are extensively used in Greek as well as in other Mediterranean languages, since similar observations have been mentioned in the use of Spanish, Portuguese and Italian. Cultural writings and conversational contexts are widely employing tenses like past continuous. It is therefore clear, that the exclusion of these suffixes makes the algorithm incomplete and inaccurate and limits its performance.

In addition to the incompletenesses of the algorithmic design, the implementation of the algorithm offers limited usability. Ntais [2008] has provided a web interface written in Javascript. Through a form, users can insert a Greek word in upper case and have its stem returned. Despite the fact that everyone can examine the algorithm, since it is embedded in the web page, it can not be directly used by any kind of application that requires stemming in Greek. Its implementation language, Javascript, is a powerful language for client side scripting in web applications. Nevertheless, it can not be used for writing a library that can be used by other applications.

### **1.7 Aims of this Thesis**

We are conducting this work in order to fully test the existing algorithm and improve it. From our initial, undocumented tests, we concluded that the existing algorithm is giving satisfactory but inaccurate and incomplete results . We are convinced that documenting and analysing the results and improving the algorithm would be a contribution not only to computer science and computational linguistics in particular but also to all these fields that Greek is used including medicine and mathematics. In addition, the algorithm can later on be used at a production level in a search engine, with the potential to give better results and a better web searching experience to users searching for documents written in Greek.

In addition to improving the existing algorithm, we will also provide a library version of the algorithm written in PHP. The reason is that PHP is currently one of the most widely used languages in the web by providing server-side scripting for building dynamic web

sites. By implementing a PHP algorithm, our aim is to provide a stemmer that can be directly used by the engine of any web application, for any kind of web search or linguistics. Other programming languages such as Javascript or even the more powerful like C and C++ lack this ability [Lerdorf and Tatroe, 2002; Flanagan, 2004]. In conclusion, our work, which will be available under an Open Source licence, will lead to a more powerful, more complete and more consistent Greek stemmer that can be directly examined, used and modified by others.

### **1.8 Research Question**

In search for solutions to the previously stated problems, the research question to be answered in this thesis can be formulated as follows:

*Up to which point the addition of more syntactical rules and exceptions improves the precision of the Ntais stemming algorithm?*

The previous algorithm by Ntais [2006] does not, deliberately, include some suffixes in an attempt to avoid errors that occur when the appropriate exception list is not also introduced with the addition of a new suffix. The creation of an exception list is a trivial but rather time consuming process. We need to identify whether it is feasible or not to create extensive and complete lists of exceptions for new rules at this point, where the existing stemmer is already producing somewhat satisfactory results and covers most of the cases.

### **1.9 Methodology**

One of the initial aims of this work was to test the original stemmer in combination with a search engine, and Google's search engine was a candidate. A web interface that would feed Google with modified, stemmed queries and unmodified ones could be easily built. The results of both modified and unmodified queries could then be compared. Unfortunately, not only the application of a stemmer in a web search engine is beyond the time limitations of this thesis, but it is also unclear whether Google is already utilizing any kind of stemming techniques for Greek. Furthermore, in a previous web search engine evaluation [Lazarinis, 2005], it is pointed out that Google returns a different number of results for different variations of the word “Athens” ( Αθήνα: Athens, Αθήνας: of Athens, Αθηνών : of (the city of) Athens). The difference in results can only imply that no stemming is used. Despite that, there are reports that some form

of stemming is being conducted [Google, 2003] although it is unclear how extensive. In addition Paice [1994] suggests that evaluating a stemmer solely in terms of IR is incomplete since IR is only one field that stemming can be applied and “...*gives no insight into the specific causes of errors*”. Because of all these reasons and due to time limitations we decided not to test our implementations of both the existing algorithm and our improved version with a search engine.

During our research, we will modify Ntais' algorithm to use more grammatical rules, exceptions and stop words. We will improve the algorithm in a constructive and extended manner. More additions will be implemented incrementally, after testing all previous improvements each time. The task of introducing more grammatical rules is challenging since it utilises techniques and requires knowledge from two domains, computer science and linguistics. In order to evaluate Ntais' and our revised algorithm, we will execute both of them in batch mode against a collection of more than half a million Greek words. Both algorithms will stem the input words from the text, and will form groups of words that have the same stem. Our purpose is to manually check whether all the words reduced by the algorithm to the same stem also share the same semantic meaning. The two algorithms are evaluated separately and the results will be compared.

## **1.10 Overview of Thesis' Contents**

In Chapter 2 we will provide a short overview of the history of the Greek language and some of the Greek grammar features and peculiarities that Greek words have during conjugation. Chapter 3 introduces some stemmer performance metrics that will be used during our evaluation in order to compare the output of the original stemmer and our modified version. The design of the existing algorithm and an extensive list of its rules are given in Chapter 4. This chapter also deals with some alternative approaches which we also considered for the re-design of our algorithm and the reasons for which we rejected them, based on their suitability to the application domain.

In Chapter 5, the evaluation of the existing algorithm is described and in Chapter 6 we describe the improvements that we decided to incorporate in our new algorithm along with the new set of rules and exception lists. The final evaluation of the algorithm and a comparison against its predecessor is given in Chapter 7.

## 2. The Greek Language

### 2.1 The History of the Greek Language

The earliest traces of written Greek can be found in more than 4400 clay tablets of the Linear B script, which was deciphered during 1951 to 1953 by the architect M. Ventris in England. This form of writing was used from 1600 to 1100 BC, and it is considered as the “earliest European script we can understand” [Robinson, 1995]. For the next 300 years, a period regarded as “the Dark Age” of illiteracy in Greece [Robinson, 1995], no traits of the Greek language have been discovered. During this period, the Homeric Greeks gave their position to the classical Greeks. The classic period of Ancient Greece (500-323 BC) coincides with the emergence of a new alphabet borrowed from the Phoenicians. Although it is debatable whether the Phoenicians or Greeks living in Phoenicia were the creators of this alphabet [Robinson, 1995], this alphabet is the ancestor of not only the Greek alphabet, but through the Etruscan and Latin languages, the ancestor of modern European alphabets.[Baugh & Cable, 2001]. The known fact is that the first consonant-only based Phoenician alphabet came to Greece without vowels, and the ancient Greeks added vowels to it. These added letter-characters improved greatly the communication and increased the use of the alphabet in everyday life and in writing form. After all, the enhanced with vowels new alphabet came nearer to the needs of everyday speech and it mirrored the spoken words more clearly and more accurately than its consonant based predecessor.

After many additions and simplifications through the years, the nowadays alphabet contains twenty four (24) letters. Table 4 presents the latest form of the Greek alphabet including both upper and lower case letters and their equivalent sound in the English alphabet. The third character of Sigma (ς) is only used in the end of a lower case word.

**Table 4: The Greek Alphabet. Letters and their equivalent sound in English**

Alpha	Beta	Gamma	Delta	Epsilon n	Zeta	Eta	Theta	Iota	Kapa	Lamda	Mi
A α	B β	Γ γ	Δ δ	E ε	Z ζ	H η	Θ θ	I ι	K κ	Λ λ	M μ
a	v	-	th (the)	e	z	i	th	i	k	l	m
Ni	Xi	Omikron	Pi	Ro	Sigma	Tau	Ypsilon	Phi	Chi	Psi	Omega
N ν	Ξ ξ	Ο ο	Π π	Ρ ρ	Σ σ ζ	Τ τ	Υ υ	Φ φ	Χ χ	Ψ ψ	Ω ω
n	ks	o	p	r	s	t	i	f	h	ps	o

The Ancient Greeks of the Classical period were organized in city states. Each of the main city states had its own region of influence and with that a different dialect. The differences between these dialects were minor, so Greek was considered as a common language [Triantafyllidis, 1941]. It was only after the conquests of Alexander the Great in Asia when the Athenian dialect, after borrowing words from other dialects, became the common dialect spoken from Greece and Egypt to Syria and Persia. The language continued to evolve for the next centuries, until the fall of Constantinople in 1453 and the beginning of the Ottoman era. For the next 400 years, following the closing of schools, the Greek language is kept oral, divided into local dialects. A few examples of written material from this period can be credited to Greeks living in countries in Central and Western Europe, like Romania, Austria, Russia, the Hungarian Empire and Italy. After the Greek War of Independence, which started in 1821, and the liberation in 1829, two competing varieties are found. The popular and spoken “Dimotiki”, used among the people, and the official and most resembling to Ancient Greek “Katharevousa”, used mostly by the intellectuals. Today, modern mainstream Greek is based on “Dimotiki” and is the official state language, since 1975, with simplifications in the intonation

system since 1981. In places, there are still local dialects with varying degrees of differentiation from the mainstream language.

Greek is spoken by 14 to 17 million people, officially in Greece and Cyprus and unofficially in countries like Australia, USA and Canada, where there are Greek and Cypriot communities. In addition, medical and philosophical terms are often Greek, Greek-derived or combinations of Latin and Greek words [Kurz & Kilian, 2001]. The domain of Humanities is undeniably a language world with terms and concepts that have originally been founded in the subject of Philosophy and Mathematics and expressed in Greek.

## 2.2 Stemming in Greek

The Greek language is grammatically more complex than English. It has conjugations and morphologically complex words [Mackridge, 1987]. Articles, adjectives, nouns and even first names and surnames may be in various cases (like nominative and genitive), in singular or plural form and they are differentiated according to their gender (masculine, feminine, neuter). Table 5 contains the singular and plural numbers of all cases and genders that the word “cat” that be found in Greek.

**Table 5: The noun “cat” in Greek**

Cases	Singular			Plural		
	masculine	feminine	neuter	masculine	feminine	neuter
nominative	γάτος	γάτα	γατί	γάτοι	γάτες	γατιά
genitive	γάτου	γάτας	γατιού	γάτων	γάτων	γατιών
accusative	γάτο	γάτα	γατί	γάτους	γάτες	γατιά
vocative	γάτε	γάτα	γατί	γάτοι	γάτες	γατιά

In addition to the complexity mentioned above, verbs are also heavily inflected. The



Greek language consists of present, past and future tenses with both perfective and imperfective aspects. These tenses have active and passive voices. Verbs are divided into two conjugations classes which have different endings and many times there are alternative endings for the same number and person. From the tables given in Appendix A, containing some examples about verb conjugation, it is obvious that Greek is much more complicated than English. Where in English four (4) endings are used, in Greek the distinct endings are 107, even without counting the different endings because of the moving mark phenomenon. Not only a stemmer has to deal with an enormously greater number of endings, but a somewhat perfect stemmer should be aware of how to deal with the “moving” tone mark issue.

### 3. Stemmer Performance Metrics

In order to evaluate the existing stemmer and measure its effectiveness, we will introduce some of the metrics that can be found in the previous literature.

#### 3.1 Frakes' Metrics

Frakes [2003] defines *stemmer strength* as the degree to which a stemmer changes words. These changes fall into two categories, *removal* and *recording*. *Removal* is the decrease of a word's length due to elimination of an affix, whereas *recording* is the replacing of a word's letter with another. Since the strength of a stemmer can affect the precision and recall in queries, Frakes defines a set of metrics that help to compare algorithms by having the algorithms stem the same texts and compare the results of the following metrics:

- The Mean Number of Words per Conflation Class

The mean number of words per conflation class is the average number of words that are found in each conflation class. For example if the words "engineer," "engineered," and "engineering" are stemmed to "engineer," then this conflation class size is three.

Stronger stemmers produce conflation classes with more words than lighter stemmers, from the same text.

- Index Compression Factor

The index compression factor is defined as  $ICS = \frac{n-s}{n}$ , where  $n$  is the number of words in the corpus and  $s$  is the number of stems produced by the stemmer. This metric indicates the index reduction that can be achieved through stemming. For example, if during the stemming of a corpus with 1000 words ( $n$ ) we end up with 800 stems ( $s$ ), we have eliminated 200 words which means an index compression factor of 20%. Stronger stemmers will have a larger index compression factor than lighter stemmers.

- The Number of Words and Stems that Differ

Stemmers often leave words unchanged. The reason behind this behaviour can be either the lack of an appropriate rule, a software bug in the implementation of the algorithm or a design choice from the authors. For example, a stemmer might not alter "engineer" because it is already a dictionary entry. A big ratio of unchanged words to total words can indicate poor algorithmic performance. Stronger stemmers will change words more

often than weaker stemmers.

- The median and mean modified Hamming distance

The Hamming distance between two strings of equal length is defined as the number of characters in the two strings that are different at the same position. For strings of unequal length we add the difference in length to the Hamming distance to give a modified Hamming distance function  $d$ . This measure takes into account transformations of stem endings. For example, a stemming algorithm might reduce the corpus { try, tried, trying } to the stem “tri”. The mean modified Hamming distance between the original words and the stem is  $D = (1+2+4)/3 = 2.33$  characters, and the median is 2.

### 3.2 Error Metrics

There are two clearly distinct error metrics categories concerning stemmers, *understemming* and *overstemming*.

Understemming occurs when words are not fully stemmed to their potential stem. In that case, words that share the same conceptual meaning are stemmed to different stems and assigned to a different conflation class.

In contrast, overstemming occurs when words that do not share the same conceptual meaning are reduced into the same stem and assigned to the same conflation class.

According to other evaluations [Alvares *et al.*, 2005], the most accurate way to check for understemming and overstemming errors is through human interference. Some examples for both categories are given in Table 6. The first example shows understemming where two words have different stems although they should had the same, since they both have to do with “selling”. On the other hand, the second examples demonstrates overstemming where two words with different semantics, “selling “ and “bird”, are incorrectly reduced to the same stem.

**Table 6: Stemming Errors**

<b>Understemming</b>
----------------------

<b>Word</b>	<b>Meaning</b>	<b>Stem</b>
πουλάω	(I) sell	πουλ-
πουλώντας	selling	πουλοντ-
<b>Overstemming</b>		
<b>Word</b>	<b>Meaning</b>	<b>Stem</b>
πουλάω	I sell	πουλ-
πουλί	bird	πουλ-

## 4. Algorithmic Design

### 4.1 The Existing Algorithm

The algorithm provided by Ntais deals “with each suffix individually” in a decentralized manner [Ntais, 2006]. The algorithm has 29 rules that treat 158 suffixes. Every rule is executed in an individual step and a set of suffixes is provided in order to remove the longest matching suffix. In all but the first steps, a list of exceptions is also examined and some different suffixes are added to the stem if needed in order to deal with the complexity of the Greek language. Additionally, each step may have its own exceptions. We have decided to keep this design and base our work on this.

Table 7 presents the algorithm by Ntais [2008]. In each step the rule with suffixes to be examined, along with actions to be taken and exceptions to be considered are described.

**Table 7: Ntais' Algorithm**

Step #	Rule	Action	Exceptions
1	Word ends in: <b>ΦΑΓΙΑ ΦΑΓΙΟΥ </b> <b>ΦΑΓΙΩΝ ΣΚΑΓΙΑ </b> <b>ΣΚΑΓΙΟΥ ΣΚΑΓΙΩΝ </b> <b>ΟΛΟΓΙΟΥ ΟΛΟΓΙΑ </b> <b>ΟΛΟΓΙΩΝ ΣΟΓΙΟΥ </b> <b>ΣΟΓΙΑ ΣΟΓΙΩΝ </b> <b>ΤΑΤΟΓΙΑ ΤΑΤΟΓΙΟΥ </b> <b>ΤΑΤΟΓΙΩΝ ΚΡΕΑΣ </b> <b>ΚΡΕΑΤΟΣ ΚΡΕΑΤΑ </b> <b>ΚΡΕΑΤΩΝ ΠΕΡΑΣ </b> <b>ΠΕΡΑΤΟΣ ΠΕΡΑΤΑ </b> <b>ΠΕΡΑΤΩΝ ΤΕΡΑΣ </b> <b>ΤΕΡΑΤΟΣ ΤΕΡΑΤΑ </b> <b>ΤΕΡΑΤΩΝ ΦΩΣ ΦΩΤΟΣ </b> <b>ΦΩΤΑ ΦΩΤΩΝ </b> <b>ΚΑΘΕΣΤΩΣ </b> <b>ΚΑΘΕΣΤΩΤΟΣ </b> <b>ΚΑΘΕΣΤΩΤΑ </b> <b>ΚΑΘΕΣΤΩΤΩΝ </b> <b>ΓΕΓΟΝΟΣ </b> <b>ΓΕΓΟΝΟΤΟΣ </b> <b>ΓΕΓΟΝΟΤΑ </b> <b>ΓΕΓΟΝΟΤΩΝ</b>	Replace suffix with:  <b>ΦΑ ΣΚΑ </b> <b>ΟΛΟ ΣΟ </b> <b>ΤΑΤΟ </b> <b>ΚΡΕ ΠΕΡ </b> <b>ΤΕΡ ΦΩ </b> <b>ΚΑΘΕΣΤ </b> <b>ΓΕΓΟΝ</b>	

2a	Word ends in: <b>ΑΔΕΣ ΑΔΩΝ</b>	Remove suffix / Check exceptions	If after removal the word does not end in: <b>ΟΚ ΜΑΜ ΜΑΝ ΜΠΑΜΠ </b> <b>ΠΑΤΕΡ ΓΙΑΓΙ ΝΤΑΝΤ ΚΥΡ ΘΕΙ </b> <b>ΠΕΘΕΡ</b>  Add “ <b>ΑΔ</b> ” in the end
2b	Word ends in: <b>ΕΔΕΣ ΕΔΩΝ</b>	Remove suffix / Check exceptions	If after removal the word ends in: ΟΠ ΙΠ ΕΜΠ ΥΠ ΓΗΠ  ΔΑΠ ΚΡΑΣΠ ΜΙΑ  Add “ <b>ΕΔ</b> ” in the end
2c	Word ends in: <b>ΟΥΔΕΣ ΟΥΔΩΝ</b>	Remove suffix / Check exceptions	If after removal the word ends in: <b>ΑΡΚ ΚΑΛΙΑΚ ΠΕΤΑΛ ΛΙΧ </b> <b>ΠΛΕΧ ΣΚ Σ ΦΑ ΦΡ ΒΕΛ ΛΟΥΑ </b> <b>ΧΝ ΣΠ ΤΡΑ ΦΕ</b>  Add “ <b>ΟΥΔ</b> ” in the end
2d	Word ends in: <b>ΕΩΣ ΕΩΝ</b>	Remove suffix / Check exceptions	If after removal the word is one of: <b>ΑΡΚ ΚΑΛΙΑΚ ΠΕΤΑΛ ΛΙΧ </b> <b>ΠΛΕΞ ΣΚ Σ ΦΑ ΦΡ ΒΕΛ ΛΟΥΑ </b> <b>ΧΝ ΣΠ ΤΡΑΓ ΦΕ</b>  Add “ <b>Ε</b> ” in the end
3	Word ends in: <b>ΙΑ ΙΟΥ ΙΩΝ</b>	Remove suffix / Check exceptions	If after removal the word ends in Vowel  Add “ <b>Ι</b> ” in the end
4	Word ends in: <b>ΙΚΑ ΙΚΟ ΙΚΟΥ ΙΚΩΝ</b>	Remove suffix / Check exceptions	If after removal the word ends in Vowel  OR is one of : <b>ΑΙ ΑΔ ΕΝΔ ΑΜΑΝ </b> <b>ΑΜΜΟΧΑΙ ΗΘ ΑΝΗΘ </b> <b>ΑΝΤΙΑ ΦΥΣ ΒΡΩΜ ΓΕΡ </b> <b>ΕΞΩΔ ΚΑΛΠ ΚΑΛΛΙΝ </b> <b>ΚΑΤΑΔ ΜΟΥΑ ΜΠΑΝ </b> <b>ΜΠΑΓΙΑΤ ΜΠΟΛ ΜΠΟΣ </b> <b>ΝΙΤ ΞΙΚ ΣΥΝΟΜΗΑ </b> <b>ΠΕΤΣ ΠΙΤΣ ΠΙΚΑΝΤ </b> <b>ΠΛΙΑΤΣ ΠΟΣΤΕΑΝ </b> <b>ΠΡΩΤΟΔ ΣΕΡΤ ΣΥΝΑΔ </b> <b>ΤΣΑΜ ΥΠΟΔ ΦΙΛΟΝ </b> <b>ΦΥΛΟΔ ΧΑΣ</b>  Add “ <b>ΙΚ</b> ” in the end

5a	Word is <i>ΑΓΑΜΕ</i>	Stem is <i>ΑΓΑΜ</i>	
5a	Word ends in : <i>ΑΓΑΜΕ ΗΣΑΜΕ </i> <i>ΟΥΣΑΜΕ Η*ΚΑΜΕ </i> <i>ΗΘΗΚΑΜΕ</i>	Remove suffix / Check exceptions	
5a	Word ends in: <i>ΑΜΕ</i>	Remove suffix / Check exceptions	If after removal the word is one of: <i>ΑΝΑΠ ΑΠΟΘ ΑΠΟΚ ΑΠΟΣΤ </i> <i>ΒΟΥΒ ΞΕΘ ΟΥΛ ΠΕΘ ΠΙΚΡ </i> <i>ΠΟΤ ΣΙΧ Χ</i>  Add “ <i>ΑΜ</i> ” in the end
5b	Word ends in: <i>ΑΓΑΝΕ ΗΣΑΝΕ </i> <i>ΟΥΣΑΝΕ ΙΟΝΤΑΝΕ </i> <i>ΙΟΤΑΝΕ ΙΟΥΝΤΑΝΕ </i> <i>ΟΝΤΑΝΕ ΟΤΑΝΕ </i> <i>ΟΥΝΤΑΝΕ ΗΚΑΝΕ </i> <i>ΗΘΗΚΑΝΕ</i>	Remove suffix / Check exceptions	If after removal the word one of: <i>ΤΡ ΤΣ</i>  Add “ <i>ΑΓΑΝ</i> ”

5b	<p>Word ends in:  <b>ΑΓΑΝΕ ΗΣΑΝΕ </b>  <b>ΟΥΣΑΝΕ ΙΟΝΤΑΝΕ </b>  <b>ΙΟΤΑΝΕ ΙΟΥΝΤΑΝΕ </b>  <b>ΟΝΤΑΝΕ ΟΤΑΝΕ </b>  <b>ΟΥΝΤΑΝΕ ΗΚΑΝΕ </b>  <b>ΗΘΗΚΑΝΕ</b></p>	<p>Remove  suffix /  Check  exceptions</p>	<p><input type="checkbox"/> If after removal the word  ends in Vowel without “Υ”  OR  is one of :</p> <p><b>ΒΕΤΕΡ ΒΟΥΛΚ ΒΡΑΧΜ Γ </b>  <b>ΔΡΑΔΟΥΜ Θ ΚΑΛΠΟΥΖ </b>  <b>ΚΑΣΤΕΛ ΚΟΡΜΟΡ ΛΑΟΠΑ </b>  <b>ΜΩΑΜΕΘ Μ ΜΟΥΣΟΥΑΜ Ν </b>  <b>ΟΥΑ Π ΠΕΛΕΚ ΠΑ ΠΟΛΙΣ </b>  <b>ΠΟΡΤΟΛ ΣΑΡΑΚΑΤΣ ΣΟΥΑΤ </b>  <b>ΤΣΑΡΛΑΤ ΟΡΦ ΤΣΙΓΓ ΤΣΟΠ </b>  <b>ΦΩΤΟΣΤΕΦ Χ ΨΥΧΟΠΑ ΑΓ </b>  <b>ΟΡΦ ΓΑΛ ΓΕΡ ΔΕΚ ΔΙΠΑ </b>  <b>ΑΜΕΡΙΚΑΝ ΟΥΡ ΠΙΘ ΠΟΥΡΙΤ </b>  <b>Σ ΖΩΝΤ ΙΚ ΚΑΣΤ ΚΟΠ ΛΙΧ </b>  <b>ΛΟΥΘΗΡ ΜΑΙΝΤ ΜΕΛ ΣΙΓ ΣΠ </b>  <b>ΣΤΕΓ ΤΡΑΓ ΤΣΑΓ Φ ΕΡ ΑΔΑΠ </b>  <b>ΑΘΙΓΓ ΑΜΗΧ ΑΝΙΚ ΑΝΟΡΓ </b>  <b>ΑΠΗΓ ΑΠΙΘ ΑΤΣΙΓΓ ΒΑΣ </b>  <b>ΒΑΣΚ ΒΑΘΥΓΑΛ ΒΙΟΜΗΧ </b>  <b>ΒΡΑΧΥΚ ΔΙΑΤ ΔΙΑΦ ΕΝΟΡΓ </b>  <b>ΘΥΣ ΚΑΠΝΟΒΙΟΜΗΧ </b>  <b>ΚΑΤΑΓΑΛ ΚΑΙΒ ΚΟΙΛΑΡΦ ΛΙΒ </b>  <b>ΜΕΓΛΟΒΙΟΜΗΧ </b>  <b>ΜΙΚΡΟΒΙΟΜΗΧ ΝΤΑΒ </b>  <b>ΞΗΡΟΚΑΙΒ ΟΛΙΓΟΔΑΜ </b>  <b>ΟΛΟΓΑΛ ΠΕΝΤΑΡΦ ΠΕΡΗΦ </b>  <b>ΠΕΡΙΤΡ ΠΛΑΤ ΠΟΛΥΔΑΠ </b>  <b>ΠΟΛΥΜΗΧ ΣΤΕΦ ΤΑΒ ΤΕΤ </b>  <b>ΥΠΕΡΗΦ ΥΠΟΚΟΠ </b>  <b>ΧΑΜΗΛΟΔΑΠ ΨΗΛΟΤΑΒ</b></p> <p>Add “AN” in the end</p>
5c	<p>Word ends in:  <b>ΗΣΕΤΕ</b></p>	<p>Remove  suffix /  Check  exceptions</p>	



5c	Word ends in: <b>ETE</b>	Remove suffix / Check exceptions  <input type="checkbox"/>	<input type="checkbox"/> If after removal the word ends in Vowel without “Y” OR is one of :  <input type="checkbox"/> <b>ΑΒΑΡ ΒΕΝ ΕΝΑΡ ΑΒΡ ΑΔ ΑΘ ΑΝ ΑΠΔ ΒΑΡΟΝ ΝΤΡ ΣΚ ΚΟΠ ΜΠΟΡ ΝΙΦ ΠΑΓ ΠΑΡΑΚΑΔ ΣΕΡΠ ΣΚΕΔ ΣΥΡΦ ΤΟΚ ΥΔ ΕΜ ΘΑΡΡ Θ</b> OR ends in : <b>ΟΔ ΑΙΡ ΦΟΡ ΤΑΘ ΔΙΑΘ ΣΧ ΕΝΔ ΕΥΡ ΤΙΘ ΥΠΕΡΘ ΡΑΘ ΕΝΘ ΡΟΘ ΣΘ ΠΥΡ ΑΙΝ ΣΥΝΔ ΣΥΝ ΣΥΝΘ ΧΩΡ ΠΟΝ ΒΡ ΚΑΘ ΕΥΘ ΕΚΘ ΝΕΤ ΡΟΝ ΑΡΚ ΒΑΡ ΒΟΔ ΩΦΕΔ</b>  Add “ <b>ET</b> ” in the end
5d	Word ends in: <b>ΟΝΤΑΣ ΩΝΤΑΣ</b>	Remove suffix / Check exceptions	If after removal the word is: <b>ΑΡΧ</b>  add “ <b>ONT</b> ” in the end OR If after removal the word is: <b>ΚΡΕ</b>  add “ <b>ΩNT</b> ” in the end
5e	Word ends in: <b>ΟΜΑΣΤΕ ΙΟΜΑΣΤΕ</b>	Remove suffix / Check exceptions	If after removal the word is: <b>ΟΝ</b>  add “ <b>ΟΜΑΣΤ</b> ” in the end
5f	Word ends in: <b>ΙΕΣΤΕ</b>	Remove suffix / Check exceptions	If after removal the word is one of: <b>Π ΑΠ ΣΥΜΠ ΑΣΥΜΠ ΑΚΑΤΑΠ ΑΜΕΤΑΜΦ</b>  Add “ <b>ΙΕΣΤ</b> ” in the end
5f	Word ends in: <b>ΕΣΤΕ</b>	Remove suffix / Check exceptions	If after removal the word is one of: <b>ΑΔ ΑΡ ΕΚΤΕΔ Ζ Μ Ξ ΠΑΡΑΚΑΔ ΑΡ ΠΡΟ ΝΙΣ</b>  Add “ <b>ΕΣΤ</b> ” in the end
5g	Word ends in: <b>ΗΘΗΚΑ ΗΘΗΚΕΣ ΗΘΗΚΕ</b>	Remove suffix / Check exceptions	

5g	Word ends in: <b>ΗΚΑ ΗΚΕΣ ΗΚΕ</b>	Remove suffix / Check exceptions  <input type="checkbox"/>	<input type="checkbox"/> If after removal the word is one of:  <b>ΔΙΑΘ Θ ΠΑΡΑΚΑΤΑΘ ΠΡΟΣΘ ΣΥΝΘ</b>  OR ends in:  <b>ΣΚΩΛ ΣΚΟΥΛ ΝΑΡΘ ΣΦ ΟΘ ΠΙΘ</b>  <input type="checkbox"/> Add “ <b>ΗΚ</b> ” in the end
5h	Word ends in: <b>ΟΥΣΑ ΟΥΣΕΣ ΟΥΣΕ</b>	Remove suffix / Check exceptions  <input type="checkbox"/>	<input type="checkbox"/> If after removal the word is one of:  <b>ΦΑΡΜΑΚ ΧΑΛ ΑΓΚ ΑΝΑΡΡ ΒΡΟΜ ΕΚΛΙΠ ΑΜΠΙΛ ΛΕΧ Μ ΠΑΤΡ Λ ΜΕΛ ΜΕΣΑΖ ΥΠΟΤΕΙΝ ΑΜ ΑΙΘ ΑΝΗΚ ΔΕΣΠΟΖ ΕΝΔΙΑΦΕΡ ΔΕ ΔΕΥΤΕΡΕΥ ΚΑΘΑΡΕΥ ΠΛΕ ΤΣΑ</b>  OR ends in:  <b>ΠΟΔΑΡ ΒΛΕΠ ΠΑΝΤΑΧ ΦΡΥΔ ΜΑΝΤΙΑ ΜΑΛΛ ΚΥΜΑΤ ΛΑΧ ΛΗΓ ΦΑΓ ΟΜ ΠΡΩΤ</b>  <input type="checkbox"/> Add “ <b>ΟΥΣ</b> ” in the end

5i	Word ends in: <b>ΑΓΑ ΑΓΕΣ ΑΓΕ</b>	Remove suffix / Check exceptions	<input type="checkbox"/> If after removal the word is one of:  <b>ΑΒΑΣΤ ΠΟΛΥΦ ΑΔΗΦ ΠΑΜΦ Ρ ΑΣΠ ΑΦ ΑΜΑΛ ΑΜΑΛΛΙ ΑΝΥΣΤ ΑΠΕΡ ΑΣΠΑΡ ΑΧΑΡ ΔΕΡΒΕΝ ΔΡΟΣΟΠ ΞΕΦ ΝΕΟΠ ΝΟΜΟΤ ΟΛΟΠ ΟΜΟΤ ΠΡΟΣΤ ΠΡΟΣΩΠΟΠ ΣΥΜΠ ΣΥΝΤ Τ ΥΠΟΤ ΧΑΡ ΑΕΙΠ ΑΙΜΟΣΤ ΑΝΥΠ ΑΠΟΤ ΑΡΤΙΠ ΔΙΑΤ ΕΝ ΕΠΙΤ ΚΡΟΚΑΛΟΠ ΣΙΔΗΡΟΠ Α ΝΑΥ ΟΥΛΑΜ ΟΥΡ Π ΤΡ Μ</b>  <input type="checkbox"/>  OR ends in:  <b>ΟΦ ΠΕΛ ΧΟΡΤ ΛΛ ΣΦ ΡΠ ΦΡ ΠΡ ΛΟΧ ΣΜΗΝ</b> BUT  is not one of: <b>ΨΟΦ ΝΑΥΛΟΧ</b> AND does not end in: <b>ΚΟΛΛ</b>  Add “ΑΓ” in the end
5j	Word ends in: <b>ΗΣΕ ΗΣΟΥ ΗΣΑ</b>	Remove suffix	<input type="checkbox"/> If after removal the word is one of: <input type="checkbox"/> <b>Ν ΧΕΡΣΟΝ ΔΩΔΕΚΑΝ ΕΡΗΜΟΝ ΜΕΓΑΛΟΝ ΕΠΤΑΝ</b>  <input type="checkbox"/> <input type="checkbox"/> Add “ΗΣ” in the end
5k	Word ends in: <b>ΗΣΤΕ</b>	Remove suffix / Check exceptions	<input type="checkbox"/> If after removal the word is one of: <input type="checkbox"/> <b>ΑΣΒ ΣΒ ΑΧΡ ΧΡ ΑΠΛ ΑΕΙΜΝ ΔΥΣΧΡ ΕΥΧΡ ΚΟΙΝΟΧΡ ΠΑΛΙΜΨ</b>  <input type="checkbox"/> <input type="checkbox"/> Add “ΗΣΤ” in the end
5l	Word ends in: <b>ΟΥΝΕ ΗΣΟΥΝΕ ΗΘΟΥΝΕ</b>	Remove suffix / Check exceptions	<input type="checkbox"/> If after removal the word is one of: <input type="checkbox"/> <b>Ν Ρ ΣΠΠ ΣΤΡΑΒΟΜΟΥΤΣ ΚΑΚΟΜΟΥΤΣ ΕΞΩΝ</b>  <input type="checkbox"/> Add “ΟΥΝ” in the end

51	Word ends in: <b>ΟΥΜΕ ΗΣΟΥΜΕ  ΗΘΟΥΜΕ</b>	Remove <input type="checkbox"/> suffix / <input type="checkbox"/> Check <input type="checkbox"/> exceptions <input type="checkbox"/>	If after removal the word is one of: <b>ΠΑΡΑΣΟΥΣ Φ Χ ΩΡΙΟΠΛ ΑΖ  ΑΛΛΙΟΣΟΥΣ ΑΣΟΥΣ</b> Add “ <b>ΟΥΜ</b> ” in the end
6	Word ends in: <b>ΜΑΤΑ ΜΑΤΩΝ ΜΑΤΟΣ</b>	Remove <input type="checkbox"/> suffix / <input type="checkbox"/> Check <input type="checkbox"/> exceptions	Always Add “ <b>ΜΑ</b> ” in the end
6	Word ends in: <b>Α ΑΙΓΑΤΕ ΑΙΓΑΝ ΑΕΙ ΑΜΑΙ ΑΝ ΑΣ ΑΣΑΙ ΑΤΑΙ ΑΩ Ε ΕΙ ΕΙΣ  ΕΙΤΕ ΕΞΑΙ ΕΞ ΕΤΑΙ Ι ΙΕΜΑΙ ΙΕΜΑΣΤΕ ΙΕΤΑΙ ΙΕΣΑΙ  ΙΕΣΑΣΤΕ ΙΟΜΑΣΤΑΝ ΙΟΜΟΥΝ ΙΟΜΟΥΝΑ ΙΟΝΤΑΝ  ΙΟΝΤΟΥΣΑΝ ΙΟΣΑΣΤΑΝ ΙΟΣΑΣΤΕ ΙΟΣΟΥΝ ΙΟΣΟΥΝΑ  ΙΟΤΑΝ ΙΟΥΜΑ ΙΟΥΜΑΣΤΕ ΙΟΥΝΤΑΙ ΙΟΥΝΤΑΝ Η ΗΔΕΣ  ΗΔΩΝ ΗΘΕΙ ΗΘΕΙΣ ΗΘΕΙΤΕ ΗΘΗΚΑΤΕ ΗΘΗΚΑΝ  ΗΘΟΥΝ ΗΘΩ ΗΚΑΤΕ ΗΚΑΝ ΗΣ ΗΣΑΝ ΗΣΑΤΕ ΗΣΕΙ  ΗΣΕΣ ΗΣΟΥΝ ΗΣΩ Ο ΟΙ ΟΜΑΙ ΟΜΑΣΤΑΝ ΟΜΟΥΝ  ΟΜΟΥΝΑ ΟΝΤΑΙ ΟΝΤΑΝ ΟΝΤΟΥΣΑΝ ΟΣ ΟΣΑΣΤΑΝ  ΟΣΑΣΤΕ ΟΣΟΥΝ ΟΣΟΥΝΑ ΟΤΑΝ ΟΥ ΟΥΜΑΙ ΟΥΜΑΣΤΕ  ΟΥΝ ΟΥΝΤΑΙ ΟΥΝΤΑΝ ΟΥΣ ΟΥΣΑΝ ΟΥΣΑΤΕ Υ ΥΣ Ω ΩΝ</b>		Remove suffix
7	Word ends in: <b>ΕΣΤΕΡ ΕΣΤΑΤ ΟΤΕΡ  ΟΤΑΤ ΥΤΕΡ ΥΤΑΤ ΩΤΕΡ  ΩΤΑΤ</b>	Remove <input type="checkbox"/> suffix	

## 4.2 Rejected Designs

During the evaluation of the existing work we rejected some alternative to Ntais' algorithm designs for our stemmer.

### 4.2.1 Lovins' Design

As mentioned in Section 1.4, Lovins' algorithm uses two steps in order to remove suffixes from words. The algorithm is used for stemming in English. As a design, it offers more simplicity than the algorithm provided by Ntais, which is executed in 29 steps. During our evaluation, we implemented a design similar to Lovins' algorithm using the same list of suffixes that Ntais used. The reason was that we observed Ntais algorithm's behaviour and it removes suffixes in one or two steps on average, despite the fact that it always executes on 29 steps. Unfortunately, our implementation of a Lovins-like algorithm didn't offer any improvement. Ntais' algorithm uses an exception list after each step in order to deal with the richness and irregularities of Greek. Even if 29 steps may be regarded as poor design, since the rest of the algorithms mentioned in Table 3

execute in two to five steps, tracking and matching endings with exceptions is easily conducted. Thus, the algorithm can be studied and improved easily. Despite of the fact that eventually in our algorithm we kept Ntais' design and we even added more steps, we made sure that the algorithm is not making unnecessary checks by returning the correct stem right after matching all possible suffixes.

#### **4.2.2 Context-Free Grammar Design**

Another alternative design that we examined was a more theoretical and sophisticated one based on the theory of context-free grammars. The main idea behind the theory is that all natural languages are based on elements, which in turn can be based into other elements recursively [Lewis and Papadimitriou, 1998]. Parse trees can be used and by applying rewrite grammatical rules the elements of sentences can be constructed and, in the case of a stemmer, de-constructed and analysed.

This approach is similar to corpus-based stemming, where each word is examined and assigned to a conflation class not only according to the grammatical meaning it holds but also to its semantic meaning it contains in the text [Xu and Croft, 2000]. Our hypothesis was that an algorithm based in these principles could be trained to recognize words and build more complex conflation classes. In that way, when a word was given, the stemmer could return a list of alternative words that have already been extracted from texts and have the same meaning. Nevertheless, this approach is beyond stemming even though it could, theoretically, offer better results in search queries. Another prohibitive reason was that although there is literature available about this subject, there is no evidence, to our knowledge, of any production or working systems and stemmer implementation using this approach, at least publicly available.

#### **4.2.3 Dictionary Based Design**

A rule based stemming algorithm has to be aware of exceptions. By adding more rules to an algorithm, a researcher must also add more exceptions, in order to deal with the complexity of a natural language. One may argue that a stemmer with many exceptions resembles a dictionary stemmer, since it has embedded in its design lists of words. A dictionary stemmer could be built by using a rule-based stemmer for creating an initial list of stems, and then manually checking the stems created. We rejected that design, since our personal interest was not only to create a solution for Greek stemming but also to study the behaviour of the language.

#### **4.2.4 Krovetz's Experimental Algorithm**

Another possible use of a dictionary for our algorithmic design could be using a dictionary to check whether a rule produced a word that is a dictionary entry. Krovetz [1993] experimented with this design. In his algorithm, prior to the execution of each rule, the word is looked up in a dictionary. If the word is an entry in the dictionary, the algorithm returns it. In the opposite case, the algorithm proceeds in evaluating more rules. When the word is modified by a rule, the resulting word is again looked up in the dictionary. This pattern continues until the remaining word is an entry or no more rules can be applied.

Krovetz experimented with this design, in an attempt to deal with the aggressiveness of Porter's algorithm and ended up with a “weaker “stemmer. For example, this experimental algorithm would stem “generalisations” to “generalisation” and not to “general”. This would provide optimal results in IR since “general” is a word with multiple entries in a dictionary. One of these entries has the same semantic meaning with “generalisations” while the other deals with military ranks. A search engine using Porter's algorithm would probably retrieve documents of both categories while trying to serve a query with the term “generalisation”.

We decided not to follow this design for three reasons. Our main objection in following this design is the fact that we are conducting this work having both linguistics and IR in mind. A stemming algorithm with similar design could possibly produce better results in IR, but it would be an incomplete algorithm. In addition, this design would require many modifications in the suffix and exception lists. Krovetz's design is manipulating suffixes in a piece by piece fashion. In contrast, Ntais' algorithm is removing suffixes in one or two steps. By following Krovetz's approach, the suffix and exception list would have to be created from the beginning, a task that would require an enormous amount of effort. Finally, Krovetz reported [1993] that in many cases this design offered poor results compared to the original algorithm by Porter, even in IR uses.

#### **4.3 Evaluation of Ntais' Algorithm**

In order to evaluate the algorithm developed by George Ntais, we first ported the algorithm in PHP. Furthermore, we implemented a set of helper applications that will be using directly the algorithm and will keep statistics about the stems returned. The operating system used for the evaluation was Gentoo Linux but because of the

portability of PHP and our style of coding the source code is portable and can be used in any platform that PHP is ported to.

Our evaluation begun by executing our port of Ntais' stemmer against a list of Greek words. We set the application in a manner that words with common stems will be grouped into conflation classes and then the output will be directed into a text file. This text file was examined manually for understemming and overstemming errors. In addition to that, we used modified Hamming Distances in order to find similar stems. From our observations, we concluded that two stems with a modified Hamming Distance of four or less can be possibly merged into one, indicating an understemming or overstemming error of the stemmer that generated them. An example is the comparison of the “BAΔIZ” and “BAΔIZAT” stems, erroneously generated by the original stemmer instead of one. The first is a stem for 23 words that have to do with the verb “to walk” while the latter is the stem generated for the word “BAΔIZATE” which is the third person plural in the Past Continuous tense of the same verb. The modified Hamming Distance between them is 2. We implemented a small application that outputs pairs of stems generated by the stemmer with a modified Hamming Distance between them of 3 or less. The pairs generated were potential understemming or overstemming errors and candidates for merge into one conflation class. After thorough examination of the lists generated we begun with the modification and the improvement of the algorithm.

## 5. Our Improvements

### 5.1 Introduction of Stop-Word Elimination

*Stop-word removal* is one of the most commonly used techniques in IR [Baeza-Yates and Ribeiro-Neto, 1999; Risvik *et al.*, 2003; Lazarinis, 2007]. We use stop-word elimination in order to improve the performance of the stemming algorithm. The stop-word list mainly contains words of length of at most four letters. These words are mainly articles, adverbs and conjunctions that can not be conjugated or stemmed. In addition, some common words that can be found in Greek texts, like initials, are also added to this list. In our modified algorithm, stop-word elimination is the first step of execution, a step that not only produces better results but also improves the running time of the algorithm. These words tend to deceive the stemmer and as a result, non existing stems of minimal length are created. The initial algorithm by Ntais solved this problem by processing only words of four letters or more. Although this approach left just a few words of 3 that could be stemmed unprocessed, we decided to add a stop-word list of more than 500 words, in order to increase precision.

### 5.2 Addition of More Grammatical Rules

In the course of our evaluation, and in order to conduct our algorithm testing and comparison, we created a set of helper applications that directly use our implementations of both Ntais' and our modified algorithm. One of these applications uses as input a list of words and creates conflation classes according to the stems returned by the stemmers. These classes were manually checked for overstemming and understemming errors in a manner similar to previous literature [Alvares *et al.*, 2005; Ntais, 2006]. According to the results, more suffixes were added in order to deal with understemming. As Ntais [2006] had pointed out, the introduction of more rules for additional suffixes raises stemming errors due to overstemming. In order to deal with this situation, we additionally added more exceptions in order to deal with overstemming and keep precision at an acceptable level. Our efforts concentrated mainly on the the addition of rules that deal with past tenses as they play a great role and can be often found in Greek texts.



### 5.3 Introduction of Lower Case Letters

The initial stemming algorithm of Ntais only accepts as input words in upper case letters, as we mentioned in Section 1.6. Our algorithm is capable of handling words given in any case, upper, lower or combinations of both. The main body of the algorithm remains unchanged and all rules are still in capital letters. Before the evaluation of the rules, each letter in the given word is converted into upper case. The algorithm stores the case of each letter individually in a different variable. The algorithm examines all rules and creates the stem in upper case. Before returning the stem of the given word, a final alteration of the stem occurs as the algorithm consults the case of each letter on the stem, and alters the case of a letter if needed.

The initial algorithm was only accepting words in upper case in order to deal with the “moving” tone mark. Although in our implementation the problem of the “moving” tone-mark still remains, we decided to treat both upper case and lower case words. The complexity of the Greek language and the time limitations of this thesis prohibit any serious attempt to solve this problem. Moreover, in a hypothetical situation in which this problem was solved, the improvement in precision would be minimal. Finally we believe that since a stem is not a real word, but just a linguistic unit, returning stems in Greek without tone marks is not an important issue.

### 5.4 The Revised Algorithm

After careful examination of the output of the original stemmer, we tried to incorporate as many modifications as possible. The original algorithm has some inaccuracies but, searching for omissions and errors can be compared with searching for a needle in a haystack. Moreover, an addition of a rule that corrects some errors may create other errors, unless an appropriate exception list is also created. Nevertheless, we added more rules in order to correct wrong patterns that kept appearing in the output. One striking example was the omission of any rules for suffixes that appear in Past Continuous (*IZA, IZES, IZE, IZAME, IZATE IZAN*) and past tenses in general. This detailed work appears in Table 8.

Table 8: Our revised algorithm (additions and modifications highlighted)

Step #	Rule	Action	Exceptions
UL1	Word contains letters in lower case	Convert letters in upper case Store their position in the word	
SWR	Word is one of the Stop Word List (Appendix C)	Return word unchanged	
1	Word ends in: <i>ΦΑΓΙΑ ΦΑΓΙΟΥ  ΦΑΓΙΩΝ ΣΚΑΓΙΑ  ΣΚΑΓΙΟΥ ΣΚΑΓΙΩΝ  ΟΛΟΓΙΟΥ ΟΛΟΓΙΑ  ΟΛΟΓΙΩΝ ΣΟΓΙΟΥ  ΣΟΓΙΑ ΣΟΓΙΩΝ  ΤΑΤΟΓΙΑ ΤΑΤΟΓΙΟΥ  ΤΑΤΟΓΙΩΝ ΚΡΕΑΣ  ΚΡΕΑΤΟΣ ΚΡΕΑΤΑ  ΚΡΕΑΤΩΝ ΠΕΡΑΣ  ΠΕΡΑΤΟΣ ΠΕΡΑΤΗ  ΠΕΡΑΤΑ ΠΕΡΑΤΩΝ  ΤΕΡΑΣ ΤΕΡΑΤΟΣ  ΤΕΡΑΤΑ ΤΕΡΑΤΩΝ ΦΩΣ  ΦΩΤΟΣ ΦΩΤΑ ΦΩΤΩΝ  ΚΑΘΕΣΤΩΣ  ΚΑΘΕΣΤΩΤΟΣ  ΚΑΘΕΣΤΩΤΑ  ΚΑΘΕΣΤΩΤΩΝ  ΓΕΓΟΝΟΣ  ΓΕΓΟΝΟΤΟΣ  ΓΕΓΟΝΟΤΑ  ΓΕΓΟΝΟΤΩΝ</i>	Replace suffix with: <i>ΦΑ ΣΚΑ  ΟΛΟ ΣΟ  ΤΑΤΟ  ΚΡΕ ΠΕΡ  ΤΕΡ ΦΩ  ΚΑΘΕΣΤ  ΓΕΓΟΝ</i>	

S1	<p>Word ends in:</p> <p><b>IZA IZEΣ IZE IZAME </b>  <b>IZATE IZAN IZANE </b>  <b>IZΩ IZEΙΣ IZEΙ </b>  <b>IZΟΥΜΕ IZΕΤΕ ΙΖΟΥΝ </b>  <b>ΙΖΟΥΝΕ</b></p>	<p>Remove suffix / Check exceptions / Exit</p>	<p>If after removal the word ends in:</p> <p><b>ΑΝΑΜΠΑ ΕΜΠΑ ΕΠΑ </b>  <b>ΞΑΝΑΠΑ ΠΑ ΠΕΡΙΠΑ ΑΘΡΟ </b>  <b>ΣΥΝΑΘΡΟ ΔΑΝΕ</b></p> <p>Add “I” in the end</p> <p>If after removal the word ends in:</p> <p><b>ΜΑΡΚ ΚΟΡΝ ΑΜΠΑΡ ΑΡΡ </b>  <b>ΒΑΘΥΡΙ ΒΑΡΚ Β ΒΟΛΒΟΡ ΓΚΡ </b>  <b>ΓΛΥΚΟΡ ΓΛΥΚΥΡ ΙΜΠ Α ΛΟΥ </b>  <b>ΜΑΡ Μ ΠΡ ΜΠΡ ΠΟΛΥΡ ΠΡ </b>  <b>ΠΙΠΕΡΟΡ</b></p> <p>Add “IZ” in the end</p>
S2	<p>Word ends in:</p> <p><b>ΩΘΗΚΑ ΩΘΗΚΕΣ </b>  <b>ΩΘΗΚΕ ΩΘΗΚΑΜΕ </b>  <b>ΩΘΗΚΑΤΕ ΩΘΗΚΑΝ </b>  <b>ΩΘΗΚΑΝΕ</b></p>	<p>Remove suffix / Check exceptions / Exit</p>	<p>If after removal the word is:</p> <p><b>ΑΙ ΒΙ ΕΝ ΥΨ ΛΙ ΖΩ Σ Χ</b></p> <p>Add “ΩΝ” in the end</p>
S3	<p>Word ends in:</p> <p><b>ΙΣΑ ΙΣΕΣ ΙΣΕ ΙΣΑΜΕ </b>  <b>ΙΣΑΤΕ ΙΣΑΝ ΙΣΑΝΕ</b></p>	<p>Remove suffix / Check exceptions / Exit</p>	<p>If word is:</p> <p><b>ΙΣΑ</b></p> <p>Stem is “ΙΣ”</p> <p>If after removal the word is:</p> <p><b>ΑΝΑΜΠΑ ΑΘΡΟ ΕΜΠΑ ΕΣΕ </b>  <b>ΕΣΩΚΛΕ ΕΠΑ ΞΑΝΑΠΑ ΕΠΕ </b>  <b>ΠΕΡΙΠΑ ΑΘΡΟ ΣΥΝΑΘΡΟ </b>  <b>ΔΑΝΕ ΚΛΕ ΧΑΡΤΟΠΑ ΕΞΑΡΧΑ </b>  <b>ΜΕΤΕΠΕ ΑΠΟΚΛΕ ΑΠΕΚΛΕ </b>  <b>ΕΚΛΕ ΠΕ ΠΕΡΙΠΑ</b></p> <p>Add “I” in the end</p> <p>If after removal the word is:</p> <p><b>ΑΝ ΑΦ ΓΕ ΓΙΓΑΝΤΟΑΦ ΓΚΕ </b>  <b>ΔΗΜΟΚΡΑΤ ΚΟΜ ΓΚ Μ Π </b>  <b>ΠΟΥΚΑΜ ΟΛΟ ΛΑΡ</b></p> <p>Add “ΙΣ” in the end</p>

S4	Word ends in: <b>ΙΣΩ ΙΣΕΙΣ ΙΣΕΙ </b> <b>ΙΣΟΥΜΕ ΙΣΕΤΕ ΙΣΟΥΝ </b> <b>ΙΣΟΥΝΕ</b>	Remove suffix / Check exceptions / Exit	If after removal the word is: <b>ΑΝΑΜΠΑ ΑΘΡΟ ΕΜΠΑ ΕΣΕ </b> <b>ΕΣΩΚΛΕ ΕΠΑ ΞΑΝΑΠΑ ΕΠΕ </b> <b>ΠΕΡΙΠΑ ΑΘΡΟ ΣΥΝΑΘΡΟ </b> <b>ΔΑΝΕ ΚΛΕ ΧΑΡΤΟΠΑ ΕΞΑΡΧΑ </b> <b>ΜΕΤΕΠΕ ΑΠΟΚΛΕ ΑΠΕΚΛΕ </b> <b>ΕΚΛΕ ΠΕ ΠΕΡΙΠΑ</b>  Add “I” in the end
S5	Word ends in: <b>ΙΣΤΟΣ ΙΣΤΟΥ ΙΣΤΟ </b> <b>ΙΣΤΕ ΙΣΤΟΙ ΙΣΤΩΝ </b> <b>ΙΣΤΟΥΣ ΙΣΤΗ ΙΣΤΗΣ </b> <b>ΙΣΤΑ ΙΣΤΕΣ</b>	Remove suffix / Check exceptions / Exit	If after removal the word is: <b>Μ Π ΑΠ ΑΡ ΗΔ ΚΤ ΣΚ ΣΧ ΥΨ </b> <b>ΦΑ ΧΡ ΧΤ ΑΚΤ ΑΟΡ ΑΣΧ ΑΤΑ </b> <b>ΑΧΝ ΑΧΤ ΓΕΜ ΓΥΡ ΕΜΠ ΕΥΠ </b> <b>ΕΧΘ ΗΦΑ ΗΦΑ ΚΑΘ ΚΑΚ ΚΥΑ </b> <b>ΛΥΓ ΜΑΚ ΜΕΓ ΤΑΧ ΦΙΛ ΧΩΡ</b>  Add “ΙΣΤ” in the end  If after removal the word is: <b>ΔΑΝΕ ΣΥΝΑΘΡΟ ΚΛΕ ΣΕ </b> <b>ΕΣΩΚΛΕ ΑΣΕ ΠΛΕ</b>  Add “I” in the end

S6	<p>Word ends in:  <b>ΙΣΜΟ ΙΣΜΟΙ ΙΣΜΟΣ </b>  <b>ΙΣΜΟΥ ΙΣΜΟΥΣ ΙΣΜΩΝ</b></p>	<p>Remove  suffix /  Check  exceptions  / Exit</p>	<p>* If after removal the word is:  <b>ΑΓΝΩΣΤΙΚ ΑΤΟΜΙΚ ΓΝΩΣΤΙΚ </b>  <b>ΕΘΝΙΚ ΕΚΛΕΚΤΙΚ ΣΚΕΠΤΙΚ </b>  <b>ΤΟΠΙΚ</b></p> <p>Remove <b>“IK”</b> from the end</p> <p>*If after removal the word is:  <b>ΣΕ ΜΕΤΑΣΕ ΜΙΚΡΟΣΕ ΕΓΚΛΕ </b>  <b>ΑΠΟΚΛΕ</b></p> <p>Add <b>“ΙΣΜ”</b>in the end</p> <p>*If after removal the word is:  <b>ΔΑΝΕ ΑΝΤΙΑΔΑΝΕ</b></p> <p>Add <b>“I”</b>in the end</p> <p>*If after removal the word is:  <b>ΑΛΕΞΑΝΔΡΙΝ ΒΥΖΑΝΤΙΝ </b>  <b>ΘΕΑΤΡΙΝ</b></p> <p>Remove <b>“IN”</b> from the end</p>
S7	<p>Word ends in:  <b>ΑΡΑΚΙ ΑΡΑΚΙΑ ΟΥΔΑΚΙ </b>  <b>ΟΥΔΑΚΙΑ</b></p>	<p>Remove  suffix /  Check  exceptions  / Exit</p>	<p>*If after removal the word is:  <b>Χ Σ</b></p> <p>Add <b>“ΑΡΑΚΙ”</b>in the end</p>

S8	<p>Word ends in:  <b>AKI AKIA ITSA ITSAΣ   ITSEΣ ITΣΩN APAKI   APAKIA</b></p>	<p>Remove  suffix /  Check  exceptions  / Exit</p>	<p>*If after removal the word is:  <b>ANΘP BAMB BP KAIM KON   KOP ΛABP ΛΟΥA MEP MOYΣT   NAGKAS ΠA P PY Σ ΣK ΣOK   ΣΠΑΝ TZ ΦAPM X KAΠAK   AAIΣΦ AMBp ANΘP K ΦYA   KATPAΠ KAIM MAA ΣΛOB ΣΦ   TΣEXOΣΛOB</b></p> <p>Add “AK” in the end</p> <p>*If after removal the word is:  <b>B BAAI ΓIAN ΓA  ZHΓOYMEH   KAPA KON MAKPYH NYΦ    ΠATEP Π  ΣK TOΣ TPIPOA</b></p> <p>Add “ITΣ” in the end</p> <p>*If after removal the word end in:  <b>KOP</b></p> <p>Add “ITΣ” in the end</p>
S9	<p>Word ends in:  <b>IAIO IAIA IAIΩN</b></p>	<p>Remove  suffix /  Check  exceptions  / Exit</p>	<p>*If after removal the word is:  <b>AIΦH IP OAO ΨAA</b></p> <p>Add “IA” in the end</p> <p>*If after removal the word iend in  <b>E ΠAIXH</b></p> <p>Add “IA” in the end</p>
S10	<p>Word ends in:  <b>ISKOΣ ISKOY ISKO   ISKE</b></p>	<p>Remove  suffix /  Check  exceptions  / Exit</p>	<p>*If after removal the word is:  <b>A IB MHN P ΦPAΓK AYK OBEA</b></p> <p>Add “ISK” in the end</p>
2a	<p>Word ends in:  <b>AAEEΣ AAΩN</b></p>	<p>Remove  suffix /  Check  exceptions  / Exit</p>	<p>If after removal the word does not  end in:  <b>OK MAM MAN MΠAMΠ   ΠATEP ΓIAΓI NTANT KYR ΘEI   ΠEΘEP</b></p> <p>Add “AA” in the end</p>

2b	Word ends in: <i>ΕΛΕΣ ΕΛΩΝ</i>	Remove suffix / Check exceptions / Exit	If after removal the word ends in: <i>ΟΠ Π ΕΜΠ ΥΠ ΓΗΠ ΔΑΠ </i> <i>ΚΡΑΣΠ ΜΙΑ</i>  Add “ <i>ΕΛ</i> ” in the end
2c	Word ends in: <i>ΟΥΛΕΣ ΟΥΛΩΝ</i>	Remove suffix / Check exceptions / Exit	If after removal the word ends in: <i>ΑΡΚ ΚΑΛΙΑΚ ΠΕΤΑΛ ΛΙΧ </i> <i>ΠΛΕΧ ΣΚ Σ ΦΑ ΦΡ ΒΕΛ ΛΟΥΑ </i> <i>ΧΝ ΣΠ ΤΡΑ ΦΕ</i>  Add “ <i>ΟΥΛ</i> ” in the end
2d	Word ends in: <i>ΕΩΣ ΕΩΝ</i>	Remove suffix / Check exceptions / Exit	If after removal the word is one of: <i>ΑΡΚ ΚΑΛΙΑΚ ΠΕΤΑΛ ΛΙΧ </i> <i>ΠΛΕΞ ΣΚ Σ ΦΑ ΦΡ ΒΕΛ ΛΟΥΑ </i> <i>ΧΝ ΣΠ ΤΡΑΓ ΦΕ</i>  Add “ <i>Ε</i> ” in the end
3	Word ends in: <i>ΙΑ ΙΟΥ ΙΩΝ</i>	Remove suffix	If after removal the word ends in Vowel  Add “ <i>Γ</i> ” in the end
4	Word ends in: <i>ΙΚΑ ΙΚΟ ΙΚΟΥ ΙΚΩΝ</i>	Remove suffix / Check exceptions / Exit	If after removal the word ends in Vowel  OR one of : <i>ΑΛ ΑΔ ΕΝΔ ΑΜΑΝ ΑΜΜΟΧΑΛ </i> <i>ΗΘ ΑΝΗΘ ΑΝΤΙΑ ΦΥΣ ΒΡΩΜ </i> <i>ΓΕΡ ΕΞΩΔ ΚΑΛΠ ΚΑΛΛΙΝ </i> <i>ΚΑΤΑΔ ΜΟΥΑ ΜΙΑΝ ΜΠΑΓΙΑΤ </i> <i>ΜΠΟΑ ΜΠΟΣ ΝΙΤ ΞΙΚ </i> <i>ΣΥΝΟΜΗΛ ΠΕΤΣ ΠΙΤΣ </i> <i>ΠΙΚΑΝΤ ΠΛΙΑΤΣ ΠΟΣΤΕΑΝ </i> <i>ΠΡΩΤΟΔ ΣΕΡΤ ΣΥΝΑΔ ΤΣΑΜ </i> <i>ΥΠΟΔ ΦΙΛΙΟΝ ΦΥΛΟΔ ΧΑΣ</i>  Add “ <i>ΙΚ</i> ” in the end
5a	Word is <i>ΑΓΑΜΕ</i>	Stem is <i>ΑΓΑΜ</i> / Exit	
5a	Word ends in : <i>ΑΓΑΜΕ ΗΣΑΜΕ </i> <i>ΟΥΣΑΜΕ ΗΚΑΜΕ </i> <i>ΗΘΗΚΑΜΕ</i>	Remove suffix / Check exceptions / Exit	

5a	Word ends in: <i>AME</i>	Remove suffix / Check exceptions / Exit	If after removal the word is one of: <i>ΑΝΑΠ ΑΠΟΘ ΑΠΟΚ ΑΠΟΣΤ </i> <i>ΒΟΥΒ ΞΕΘ ΟΥΙ ΠΕΘ ΠΙΚΡ </i> <i>ΠΟΤ ΣΙΧ Χ</i>  Add “ <i>AM</i> ” in the end
5b	Word ends in: <i>ΑΓΑΝΕ ΗΣΑΝΕ </i> <i>ΟΥΣΑΝΕ ΙΟΝΤΑΝΕ </i> <i>ΙΟΤΑΝΕ ΙΟΥΝΤΑΝΕ </i> <i>ΟΝΤΑΝΕ ΟΤΑΝΕ </i> <i>ΟΥΝΤΑΝΕ ΗΚΑΝΕ </i> <i>ΗΘΗΚΑΝΕ</i>	Remove suffix / Check exceptions / Exit	If after removal the word one of: <i>ΤΡ ΤΣ</i>  Add “ <i>ΑΓΑΝ</i> ”



5b	Word ends in: <b>ΑΓΑΝΕ ΗΣΑΝΕ </b> <b>ΟΥΣΑΝΕ ΙΟΝΤΑΝΕ </b> <b>ΙΟΤΑΝΕ ΙΟΥΝΤΑΝΕ </b> <b>ΟΝΤΑΝΕ ΟΤΑΝΕ </b> <b>ΟΥΝΤΑΝΕ ΗΚΑΝΕ </b> <b>ΗΘΗΚΑΝΕ</b>	Remove suffix / Check exceptions / Exit	<input type="checkbox"/> If after removal the word ends in Vowel without “Υ” OR is one of : <b>ΒΕΤΕΡ ΒΟΥΛΚ ΒΡΑΧΜ Γ </b> <b>ΔΡΑΔΟΥΜ Θ ΚΑΛΠΟΥΖ </b> <b>ΚΑΣΤΕΛ ΚΟΡΜΟΡ </b> <b>ΛΑΟΠΑ ΜΩΑΜΕΘ Μ </b> <b>ΜΟΥΣΟΥΑΜ Ν ΟΥΑ Π </b> <b>ΠΕΛΕΚ ΠΛ ΠΟΛΙΣ </b> <b>ΠΟΡΤΟΛ ΣΑΡΑΚΑΤΣ </b> <b>ΣΟΥΑΤ ΤΣΑΡΑΑΤ ΟΡΦ </b> <b>ΤΣΙΓΓ ΤΣΟΠ </b> <b>ΦΩΤΟΣΤΕΦ Χ ΨΥΧΟΠΑ </b> <b>ΑΓ ΟΡΦ ΓΑΛ ΓΕΡ ΔΕΚ </b> <b>ΔΙΠΛ ΑΜΕΡΙΚΑΝ ΟΥΡ </b> <b>ΠΙΘ ΠΟΥΡΙΤ Σ ΖΩΝΤ ΙΚ </b> <b>ΚΑΣΤ ΚΟΠ ΛΙΧ ΛΟΥΘΗΡ </b> <b>ΜΑΙΝΤ ΜΕΛ ΣΙΓ ΣΠ </b> <b>ΣΤΕΓ ΤΡΑΓ ΤΣΑΓ Φ ΕΡ </b> <b>ΑΔΑΠ ΑΘΙΓΓ ΑΜΗΧ </b> <b>ΑΝΙΚ ΑΝΟΡΓ ΑΠΗΓ </b> <b>ΑΠΙΘ ΑΤΣΙΓΓ ΒΑΣ </b> <b>ΒΑΣΚ ΒΑΘΥΓΑΛ </b> <b>ΒΙΟΜΗΧ ΒΡΑΧΥΚ ΔΙΑΤ </b> <b>ΔΙΑΦ ΕΝΟΡΓ ΘΥΣ </b> <b>ΚΑΠΝΟΒΙΟΜΗΧ </b> <b>ΚΑΤΑΓΑΛ ΚΛΙΒ </b> <b>ΚΟΙΛΑΡΦ ΛΙΒ </b> <b>ΜΕΓΛΟΒΙΟΜΗΧ </b> <b>ΜΙΚΡΟΒΙΟΜΗΧ ΝΤΑΒ </b> <b>ΞΗΡΟΚΛΙΒ ΟΛΙΓΟΔΑΜ </b> <b>ΟΛΟΓΑΛ ΠΕΝΤΑΡΦ </b> <b>ΠΕΡΗΦ ΠΕΡΙΤΡ ΠΛΑΤ </b> <b>ΠΟΛΥΔΑΠ ΠΟΛΥΜΗΧ </b> <b>ΣΤΕΦ ΤΑΒ ΤΕΤ </b> <b>ΥΠΕΡΗΦ ΥΠΟΚΟΠ </b> <b>ΧΑΜΗΛΟΔΑΠ </b> <b>ΨΗΛΟΤΑΒ</b>
5c	Word ends in: <b>ΗΣΕΤΕ</b>	Remove suffix / Check exceptions / Exit	Add “AN” in the end

5c	Word ends in: <b>ETE</b>	Remove suffix / Check exceptions / Exit	<input type="checkbox"/> If after removal the word ends in Vowel without “Y” OR is one of : <b>ΑΒΑΡ ΒΕΝ ΕΝΑΡ ΑΒΡ ΑΔ ΑΘ ΑΝ ΑΠΔ ΒΑΡΟΝ ΝΤΡ ΣΚ ΚΟΠ ΜΠΟΡ ΝΙΦ ΠΑΓ ΠΑΡΑΚΑΔ ΣΕΡΠ ΣΚΕΔ ΣΥΡΦ ΤΟΚ ΥΔ ΕΜ ΘΑΡΡ Θ</b> OR ends in : <b>ΟΔ ΑΙΡ ΦΟΡ ΤΑΘ ΔΙΑΘ ΣΧ ΕΝΔ ΕΥΡ ΤΙΘ ΥΠΕΡΘ ΡΑΘ ΕΝΘ ΡΟΘ ΣΘ ΠΥΡ ΑΙΝ ΣΥΝΔ ΣΥΝ ΣΥΝΘ ΧΩΡ ΠΟΝ ΒΡ ΚΑΘ ΕΥΘ ΕΚΘ ΝΕΤ ΡΟΝ ΑΡΚ ΒΑΡ ΒΟΔ ΩΦΕΔ</b>  Add “ <b>ET</b> ” in the end
5d	Word ends in: <b>ΟΝΤΑΣ ΩΝΤΑΣ</b>	Remove suffix / Check exceptions / Exit	If after removal the word is: <b>ΑΡΧ</b>  add “ <b>ΟΝΤ</b> ” in the end OR If after removal the word is: <b>ΚΡΕ</b>  add “ <b>ΩΝΤ</b> ” in the end
5e	Word ends in: <b>ΟΜΑΣΤΕ ΙΟΜΑΣΤΕ</b>	Remove suffix / Check exceptions / Exit	If after removal the word is: <b>ΟΝ</b>  add “ <b>ΟΜΑΣΤ</b> ” in the end
5f	Word ends in: <b>ΙΕΣΤΕ</b>	Remove suffix / Check exceptions / Exit	If after removal the word is one of: <b>Π ΑΠ ΣΥΜΠ ΑΣΥΜΠ ΑΚΑΤΑΠ ΑΜΕΤΑΜΦ</b>  Add “ <b>ΙΕΣΤ</b> ” in the end
5f	Word ends in: <b>ΕΣΤΕ</b>	Remove suffix / Check exceptions / Exit	If after removal the word is one of: <b>ΑΔ ΑΡ ΕΚΤΕΔ Ζ Μ Ξ ΠΑΡΑΚΑΔ ΑΡ ΠΡΟ ΝΙΣ</b>  Add “ <b>ΕΣΤ</b> ” in the end

5g	Word ends in: <b>ΗΘΗΚΑ ΗΘΗΚΕΣ  ΗΘΗΚΕ</b>	Remove suffix / Check exceptions / Exit	
5g	Word ends in: <b>ΗΚΑ ΗΚΕΣ ΗΚΕ</b>	Remove suffix / Check exceptions / Exit	<input type="checkbox"/> If after removal the word is one of:  <b>ΔΙΑΘ Θ ΠΑΡΑΚΑΤΑΘ ΠΡΟΣΘ  ΣΥΝΘ</b>  OR ends in:  <b>ΣΚΩΛ ΣΚΟΥΛ ΝΑΡΘ ΣΦ ΟΘ  ΠΙΘ</b>  <input type="checkbox"/> Add “ <b>ΗΚ</b> ” in the end
5h	Word ends in: <b>ΟΥΣΑ ΟΥΣΕΣ ΟΥΣΕ</b>	Remove suffix / Check exceptions / Exit	<input type="checkbox"/> If after removal the word is one of:  <b>ΦΑΡΜΑΚ ΧΑΛ ΑΓΚ ΑΝΑΡΡ  ΒΡΟΜ ΕΚΛΙΠ ΛΑΜΠΙΔ ΛΕΧ Μ  ΠΑΤ Ρ Λ ΜΕΔ ΜΕΣΑΖ  ΥΠΟΤΕΙΝ ΑΜ ΑΙΘ ΑΝΗΚ  ΔΕΣΠΟΖ ΕΝΔΙΑΦΕΡ ΔΕ  ΔΕΥΤΕΡΕΥ ΚΑΘΑΡΕΥ ΠΛΕ ΤΣΑ</b>  OR ends in:  <b>ΠΟΔΑΡ ΒΛΕΠ ΠΑΝΤΑΧ ΦΡΥΔ  ΜΑΝΤΙΑ ΜΑΛΛ ΚΥΜΑΤ ΛΑΧ  ΛΗΓ ΦΑΓ ΟΜ ΠΡΩΤ</b>  <input type="checkbox"/> Add “ <b>ΟΥΣ</b> ” in the end

5i	Word ends in: <b>ΑΓΑ ΑΓΕΣ ΑΓΕ</b>	Remove suffix / Check exceptions / Exit	<input type="checkbox"/> If after removal the word is one of:  <b>ΑΒΑΣΤ ΠΟΛΥΦ ΑΔΗΦ ΠΑΜΦ Ρ ΑΣΠ ΑΦ ΑΜΑΛ ΑΜΑΛΛΙ ΑΝΥΣΤ ΑΠΕΡ ΑΣΠΑΡ ΑΧΑΡ ΔΕΡΒΕΝ ΔΡΟΣΟΠ ΞΕΦ ΝΕΟΠ ΝΟΜΟΤ ΟΛΟΠ ΟΜΟΤ ΠΡΟΣΤ ΠΡΟΣΩΠΟΠ ΣΥΜΠ ΣΥΝΤ Τ ΥΠΟΤ ΧΑΡ ΑΕΙΠ ΑΙΜΟΣΤ ΑΝΥΠ ΑΠΟΤ ΑΡΤΙΠ ΔΙΑΤ ΕΝ ΕΠΙΤ ΚΡΟΚΑΛΟΠ ΣΙΔΗΡΟΠ Δ ΝΑΥ ΟΥΛΑΜ ΟΥΡ Π ΤΡ Μ</b>  <input type="checkbox"/>  OR ends in:  <b>ΟΦ ΠΕΛ ΧΟΡΤ ΛΛ ΣΦ ΡΠ ΦΡ ΠΡ ΛΟΧ ΣΜΗΝ</b> BUT  is not one of: <b>ΨΟΦ ΝΑΥΛΟΧ</b> AND does not end in: <b>ΚΟΛΛ</b>  Add “ΑΓ” in the end
5j	Word ends in: <b>ΗΣΕ ΗΣΟΥ ΗΣΑ</b>	Remove suffix / Check exceptions / Exit	<input type="checkbox"/> If after removal the word is one of: <input type="checkbox"/> <b>Ν ΧΕΡΣΟΝ ΔΩΔΕΚΑΝ ΕΡΗΜΟΝ ΜΕΓΑΛΟΝ ΕΠΤΑΝ</b> <input type="checkbox"/> Add “ΗΣ” in the end
5k	Word ends in: <b>ΗΣΤΕ</b>	Remove suffix / Check exceptions / Exit	<input type="checkbox"/> If after removal the word is one of: <input type="checkbox"/> <b>ΑΣΒ ΣΒ ΑΧΡ ΧΡ ΑΠΛ ΑΕΙΜΝ ΔΥΣΧΡ ΕΥΧΡ ΚΟΙΝΟΧΡ ΠΑΛΙΜΨ</b> <input type="checkbox"/> Add “ΗΣΤ” in the end
5l	Word ends in: <b>ΟΥΝΕ ΗΣΟΥΝΕ ΗΘΟΥΝΕ</b>	Remove suffix / Check exceptions / Exit	<input type="checkbox"/> If after removal the word is one of: <input type="checkbox"/> <b>Ν Ρ ΣΠΠ ΣΤΡΑΒΟΜΟΥΤΣ ΚΑΚΟΜΟΥΤΣ ΕΞΩΝ</b> <input type="checkbox"/> Add “ΟΥΝ” in the end

51	Word ends in: <b>ΟΥΜΕ ΗΣΟΥΜΕ </b> <b>ΗΘΟΥΜΕ</b>	Remove <input type="checkbox"/> suffix / <input type="checkbox"/> Check <input type="checkbox"/> exceptions <input type="checkbox"/> / Exit <input type="checkbox"/>	If after removal the word is one of: <b>ΠΑΡΑΣΟΥΣ Φ Χ ΩΡΙΟΠΛ ΑΖ </b> <b>ΑΛΛΙΟΣΟΥΣ ΑΣΟΥΣ</b> Add “ <b>ΟΥΜ</b> ” in the end
6	Word ends in: <b>ΜΑΤΑ ΜΑΤΩΝ ΜΑΤΟΣ</b>	Remove <input type="checkbox"/> suffix <input type="checkbox"/>	Always Add “ <b>ΜΑ</b> ” in the end
6	Word ends in: <b>Α ΑΓΑΤΕ ΑΓΑΝ ΑΕΙ ΑΜΑΙ ΑΝ ΑΣ ΑΣΑΙ ΑΤΑΙ ΑΩ Ε ΕΙ ΕΙΣ </b> <b>ΕΙΤΕ ΕΣΑΙ ΕΣ ΕΤΑΙ Ι ΙΕΜΑΙ ΙΕΜΑΣΤΕ ΙΕΤΑΙ ΙΕΣΑΙ </b> <b>ΙΕΣΑΣΤΕ ΙΟΜΑΣΤΑΝ ΙΟΜΟΥΝ ΙΟΜΟΥΝΑ ΙΟΝΤΑΝ </b> <b>ΙΟΝΤΟΥΣΑΝ ΙΟΣΑΣΤΑΝ ΙΟΣΑΣΤΕ ΙΟΣΟΥΝ ΙΟΣΟΥΝΑ </b> <b>ΙΟΤΑΝ ΙΟΥΜΑ ΙΟΥΜΑΣΤΕ ΙΟΥΝΤΑΙ ΙΟΥΝΤΑΝ Η ΗΔΕΣ </b> <b>ΗΛΩΝ ΗΘΕΙ ΗΘΕΙΣ ΗΘΕΙΤΕ ΗΘΗΚΑΤΕ ΗΘΗΚΑΝ ΗΘΟΥΝ </b> <b>ΗΘΩ ΗΚΑΤΕ ΗΚΑΝ ΗΣ ΗΣΑΝ ΗΣΑΤΕ ΗΣΕΙ ΗΣΕΣ ΗΣΟΥΝ </b> <b>ΗΣΩ Ο ΟΙ ΟΜΑΙ ΟΜΑΣΤΑΝ ΟΜΟΥΝ ΟΜΟΥΝΑ ΟΝΤΑΙ </b> <b>ΟΝΤΑΝ ΟΝΤΟΥΣΑΝ ΟΣ ΟΣΑΣΤΑΝ ΟΣΑΣΤΕ ΟΣΟΥΝ </b> <b>ΟΣΟΥΝΑ ΟΤΑΝ ΟΥ ΟΥΜΑΙ ΟΥΜΑΣΤΕ ΟΥΝ ΟΥΝΤΑΙ </b> <b>ΟΥΝΤΑΝ ΟΥΣ ΟΥΣΑΝ ΟΥΣΑΤΕ Υ ΥΣ Ω ΩΝ</b>		Remove suffix
7	Word ends in: <b>ΕΣΤΕΡ ΕΣΤΑΤ ΟΤΕΡ </b> <b>ΟΤΑΤ ΥΤΕΡ ΥΤΑΤ ΩΤΕΡ </b> <b>ΩΤΑΤ</b>	Remove <input type="checkbox"/> suffix	
<b>UL2</b>	Word has converted letters to upper case from step <b>UL1</b>	<b>Convert these letters back to lower case</b>	

## 6. Final Evaluation

After the implementation of all the modifications that would improve the existing algorithm, we began to evaluate our revised algorithm. We used the same word list and the same applications we created for testing the initial algorithm. By doing so, we made sure that the statistics produced can be directly comparable.

As a result of our modifications, we ended with a stronger stemmer. Although the two stemmers leave unchanged roughly the same number of words, our modified version produces fewer and bigger conflation classes by altering more letters in every word, on average. Table 9 presents the statistics gathered after executing both algorithms against a list of 574.621 Greek words.

**Table 9: Statistics: Comparison of the original and revised algorithm**

	<b>Original Stemmer by Ntais</b>	<b>Our modified Stemmer</b>
<b>Mean number of words per conflation class</b>	4,055	5,664
<b>Index compression factor</b>	75,34%	82,34
<b>Ratio of unchanged to total words</b>	2%	2%
<b>Mean modified Humming Distance</b>	2.441	2,916
<b>Median Modified Humming Distance</b>	2	2
<b>Correct Stems (for a sample of 12468 words)</b>	10885 (87,3%)	11669 (93,52%)
<b>Distribution of Stemming errors per algorithm</b>		
	<b>Original Stemmer by Ntais</b>	<b>Our modified Stemmer</b>
<b>Understemming Errors (Section 3.2)</b>	88,44%	23,67%
<b>Overstemming Errors (Section 3.2)</b>	11,56%	76,33%
<b>Number of different stems generated by the two stemmers (for a sample of 574.621 words)</b>	35885 (6,24%)	

It is worth noting that, according to our tests, the majority of the errors of the initial algorithm had to do with understemming (88,44%). Our algorithm produces more overstemming errors (76,33%) despite of the fact that the total number of errors is reduced. The two stemmers produced 35885 different stems for the same list of words. In addition the number of executions steps was increased from 29 in the original algorithm to 42. Ten of the new execution steps have to do with the 72 new stems that were added while the remaining three deal with stop-word removal and lower to upper and upper to lower case treatment.

Although the number of steps was increased by approximately 44%, the algorithm now executes 23.17 steps on average. The reason is that while the original algorithm always executes all of its 29 steps, our modified algorithms returns the correct stem and then exits earlier if the remaining rules are not going to modify the word any further. For example, if a rule that treats verb suffixes is evaluated and executed, it is certain that no rules that deal with noun suffixes will be executed later. This modification will certainly offer better running times in any implementation of our algorithm. Nevertheless, we avoid mentioning any running times and we only refer to the average number of steps executed. The reason is that the actual running time of an algorithm depends on factors like the implementation language, coding style and coding efficiency, hardware, current load of the system and operating system among others. Because of that, we believe that comparing running times of our tests between our algorithm and the initial algorithm by Ntais would be misleading.

## 7. Conclusion

Our purpose has been to evaluate and improve the existing stemming algorithm of the Greek language. After an evaluation of the results that the original algorithm provided, we incrementally improved it by adding new rules and exceptions. Overall, we managed to gain significant improvements in completeness and accuracy to the existing algorithm.

After manually checking the output of the stemmer we conclude that our new algorithm returns more correct results than its predecessor. Due to our efforts the understemming and overstemming errors are less. The new algorithm is more complete since it supports most tenses and correctly stems suffixes that were not included before, like diminutives and others.

Moreover, we offered a more usable implementation that can be used with slight modifications or even directly. Due to the implementation language that we chose, PHP, our implementation can be used by any web or non web application that may require stemming of Greek words. Our implementation is directly usable by any kind of application, web or not, linguistic or IR related that requires stemming for Greek.

To answer our research question, we conclude that the addition of more suffixes is attainable but the effort required for each additional suffix increases geometrically. The previous algorithm already deals with the majority of suffixes that can be found in the Greek grammar. Since the majority of cases is already covered, each addition to the stemmer requires a considerable amount of evaluating the stemmer's output and searching in dictionaries for possible exceptions.

Despite of the fact that we are pleased with the outcome of our efforts, there is a lot of space for improvements. Our algorithm, like its predecessor, is not dealing with the moving tone-mark issue. Any attempt to deal with issue will probably require a considerable amount of effort due to the complexity of the Greek language. Although it will not provide any great improvement in precision, a stemmer must successfully deal with this issue in order to be considered complete. Furthermore, the stemmer can be enhanced by the addition of more suffixes and exceptions. In addition to the 158 suffixes of the initial algorithm, we added rules for 72 more. Although the majority of the cases is covered, more rules can be added in order to produce a more complete stemmer. Finally, the stemmer can be more thoroughly tested. Although we used metrics



to measure our performance and compare it with its predecessor, we had to manually check the stems generated for overstemming and understemming errors. This is a task that can only be done manually, by experts of the Greek language. Although we had an enormous collection of Greek words, we were able to test only one small portion of it, as we mention in Table 9. A more complete testing will indicate more errors and can serve as a basis for further improvement.

We believe that this thesis work contributed in stemming research by continuing and extending the work done by others in the past and by offering a stemmer for the Greek language which others can use and extend even more.

## References

- [Alvares *et al*, 2005] Alvares Reinaldo Viana , Garcia Ana Cristina Bicharra, & Ferraz Inhaúma. STEMBR: A Stemming Algorithm for the Brazilian Portuguese Language. In: *Proc. of the 12th Portuguese Conference on Artificial Intelligence, EPIA 2005, Lecture Notes in Artificial Intelligence* 3808, 693-701.
- [Baeza-Yates & Ribeiro-Neto, 1999] Baeza-Yates Ricardo & Ribeiro-Neto Berthirer. *Modern Information Retrieval*. Addison Wesley, New York, 1999.
- [Baugh & Cable, 2001] Baugh Albert & Cable Thomas. *A History of the English Language*. Prentice Hall, Upper Saddle River, 2001.
- [Brants, 2003] Brants Thorsten. Natural Language Processing in Information Retrieval. In *Proceedings of the 14th Meeting of Computational Linguistics in the Netherlands* Antwerp, Belgium.
- [Carlberger *et al*, 2001] Carlberger Johan, Dalianis Hercules, Hassel Martin & Knutsson Ola. Improving Precision in Information Retrieval for Swedish using Stemming. In *Proceedings of NODALIDA 01 - 13th Nordic Conference on Computational Linguistics*, May 21-22, Uppsala, Sweden.
- [Flanagan, 2004] Flanagan David. *Javascript: The Definitive Guide*. O'Reilly, Sebastopol , CA.
- [Frakes, 2003] Frakes William B. Strength and Similarity of Affix Removal Stemming Algorithms. *ACM SIGIR Forum archive*, 37, 1 (Spring 2003), 26 – 30.
- [Google, 2003] Google Starts Auto Stemming Searches.  
[http://www.searchengineshowdown.com/blog/2003/11/google\\_starts\\_auto\\_stemming\\_se.sh](http://www.searchengineshowdown.com/blog/2003/11/google_starts_auto_stemming_se.sh) (Retrieved 04/05/2008).
- [Krovetz, 1993] Krovetz Robert. Viewing morphology as an inference process. In *Proceedings of the 16th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, 191-202.
- [Kurz & Kilian, 2001] Kurz Thorsten & Stoffel Killian. Going beyond Stemming: Creating Concept Signatures of Complex Medical Terms. *Knowledge Based Systems Journal* 15, 5 (2001), 309-313.
- [Lazarinis, 2005] Lazarinis Fotis. Do search engines understand Greek or users requests “sound Greek” to them ? In: *M. Beigbeder and W.G. Yee (eds), Open Source Web*

*Information Retrieval Workshop in conjunction with IEEE/WIC/ACM International Conference on Web Intelligence & Intelligent Agent Technology*, Compiègne, France, 19 September 2005 (IEEE, 2005) 43-6.

[Lazarinis, 2007] Lazarinis Fotis. Lemmatization and stopword elimination in Greek web searching. *Proceedings of the 2007 Euro American Conference on Telematics and information Systems EATIS '07*, ACM, New York, NY.

[Lerdorf & Tatroe, 2002] Lerdorf Rasmus & Tatroe Kevin. *Programming PHP*. O'Reilly, Sebastopol, CA.

[Lewis & Papadimitriou, 1998] Lewis R. H. & Papadimitriou. H. C. *Elements of the Theory of Computation*. Prentice Hall, Upper Saddle River, 1998.

[Lovins, 1968] Lovins JB. Development of a stemming algorithm. *Mechanical Translation and Computational Linguistics*, 11, 1 (1968). 22-31.

[Mackridge, 1987] Mackridge P. *The Modern Greek Language. A Descriptive Analysis of Standard Modern Greek*. Clarendon Press, Oxford.

[Ntais, 2006] Ntais Georgios. *Development of a stemmer for the Greek language*. Master's Thesis at Stockholm University / Royal Institute of Technology.

[Ntais, 2008] Ntais Georgios. Online Greek Stemmer. Web Interface implemented in Javascript [http://people.dsv.su.se/~hercules/greek\\_stemmer.gr.html](http://people.dsv.su.se/~hercules/greek_stemmer.gr.html) (Retrieved 15/04/2008).

[Paice, 1994] Paice D. Chris. An Evaluation Method for Stemming Algorithms. In *Proceedings of the 17th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, 42-50.

[Paice, 1996] Paice D. Chris. Method for Evaluation of Stemming Algorithms Based on Error Counting. *Journal of the American Society for Information Science*, Volume 47, 8 (August 1996), 632 – 649.

[Porter, 1980] Porter Michael. An algorithm for suffix stripping. *Program*, 14(3), 130-137.

[Risvik *et al*, 2003] Risvik Knut Magne, Mikolajewski Tomasz & Boros Peter. Query Segmentation for Web Search. In *Proceedings of the 12th International World Wide Web Conference*, 52.

[Robinson, 1995] Robinson Andrew. *The Story of Writing*. Thames & Hudson

[Triantafyllidis, 1941] Triantafyllidis Manolis. *Modern Greek Grammar*. Institute M

Triantalyllidis.

[van Rijsbergen, 1979] van Rijsbergen, C.J. *Information Retrieval*. Butterworths, London.

[Witten et al, 2007] Witten Ian H., Gori Marco & Numerico Teresa. *Web Dragons: Inside the Myths of Search Engine Technology*. Morgan Kaufmann Publishers.

[Xu & Croft, 1998] Xu Jinxi & Croft Bruce. Corpus-Based Stemming Using Cooccurrence of Word Variants. *ACM Transactions on Information Systems*, 16, 1 (January 1998), 61-81.

[Xu & Croft, 2000] Xu Jinxi & Croft Bruce. Improving the Effectiveness of Information Retrieval with Local Context Analysis. *ACM Transactions on Information Systems* 18, 1 (January 2000), 79-112.

## Appendices

### Appendix A: Verb Conjugation Classes In Greek

#### Verbs of 1<sup>st</sup> Conjugation Classes

The verb "to read" in Greek						
Active Voice / Past Tenses						
	Υπερσυντέλικος (Past Perfect)		Αόριστος (Simple Past)		Παρατατικός (Past continuous )	
<b>Singular</b>	είχα διαβάσει	(I) had studied	διάβασα	(I) studied	διάβαζα	(I) was studying
	είχες διαβάσει	(you) had studied	διάβασες	(you) studied	διάβαζες	(you) were studying
	είχε διαβάσει	(he) had studied	διάβασε	(he) studied	διάβαζε	(he) was studying
<b>Plural</b>	είχαμε διαβάσει	(we) had studied	διαβάσαμε	(we) studied	διαβάζαμε	(we) were studying
	είχατε διαβάσει	(you) had studied	διαβάσατε	(you) studied	διαβάζατε	(you) were studying
	είχαν διαβάσει	(they) had studied	διαβάσανε (διάβασαν)	(they) studied	διάβαζαν	(they) were studying

II

The verb "to read" in Greek				
Active Voice / Present Tenses				
	Ενεστώτας (Present)		Παρακείμενος (Present Perfect)	
<b>Singular</b>	διαβάζω	(I) read /am reading	έχω διαβάσει	(I) have read /have been reading
	διαβάζεις	(you) read / are reading	έχεις διαβάσει	(you) have read/have been reading
	διαβάζει	(he) reads/ is reading	έχει διαβάσει	(he) has read/has been reading
<b>Plural</b>	διαβάζουμε	(we) read/ are reading	έχουμε διαβάσει	(we) have read/have been reading
	διαβάζετε	(you) read / are reading	έχετε διαβάσει	(you) have read/have been reading
	διαβάζουν	(they) read /are reading	έχουν διαβάσει	(they) have read/have been reading

### III

The verb “to read” in Greek						
Active Voice / Future Tenses						
	Εξακολουθητικός Μέλλοντας ( Future Continuous)		Στιγμιαίος Μέλλοντας (Simple Future)		Συντελεσμένος Μέλλοντας (Future Perfect/Imperfect)	
<b>Singular</b>	θα διαβάζω	(I) will be reading	θα διαβάσω	(I) will read	θα έχω διαβάσει	(I) will have read / have been reading
	θα διαβάζεις	(you) will be reading	θα διαβάσεις	(you) will read	θα έχεις διαβάσει	(you) will have read/have been reading
	θα διαβάζει	(he) will be reading	θα διαβάσει	(he) will read	θα έχει διαβάσει	(he) will have read/have been reading
<b>Plural</b>	θα διαβάζουμε	(we )will reading	θα διαβάσουμε	(we) will read	θα έχουμε διαβάσει	(we) will have read/have been reading
	θα διαβάζετε	(you)will be reading	θα διαβάσετε	(you) will read	θα έχετε διαβάσει	(you) will have read/have been reading
	θα διαβάζουν	(they) will be reading	θα διαβάσουνε	(they) will read	θα έχουν διαβάσει	(they) will have read/have been reading

## IV

The verb “to read” in Greek						
Passive Voice / Past Tenses						
	Υπερσυντέλικος (Past Perfect)		Αόριστος (Simple Past)		Παρατατικός (Past Continuous )	
<b>Singular</b>	είχα διαβαστεί	(I) had been studied	διαβάστηκα	(I) was studied	διαβαζόμουν	(I) was being studied
	είχες διαβαστεί	(you) had been studied	διαβάστηκες	(you) were studied	διαβαζόσουν	(you) were being studied
	είχε διαβαστεί	(he) had been studied	διαβάστηκε	(he) was studied	διαβαζόταν	(he) was being studied
<b>Plural</b>	είχαμε διαβαστεί	(we) had been studied	διαβαστήκαμε	(we) were studied	διαβαζόμαστε	(we) were being studied
	είχατε διαβαστεί	(you) had been studied	διαβαστήκατε	(you) were studied	διαβαζόσατε	(you) were being studied
	είχαν διαβαστεί	(they) had been studied	διαβαστήκανε (διαβάστηκαν)	(they) were studied	διαβαζόνταν	(they) were being studied



<b>The verb “to read” in Greek</b>				
<b>Active Voice / Present Tenses</b>				
	<b>Ενεστώτας (Present)</b>		<b>Παρακείμενος (Present Perfect)</b>	
<b>Singular</b>	διαβάζομαι	(I) am read /am being read	έχω διαβαστεί	(I) have been read
	διαβάζεσαι	(you) are read / are being read	έχεις διαβαστεί	(you) have been read
	διαβάζεται	(he) is read/ is being read	έχει διαβαστεί	(he) has been read
<b>Plural</b>	διαβάζομαστε	(we) are read/ are being read	έχουμε διαβαστεί	(we) have been read
	διαβάζεστε	(you) are read / are being read	έχετε διαβαστεί	(you) have been read
	διαβάζονται	(they) are read /are being read	έχουν διαβαστεί	(they) have been read

## VI

The verb “to read” in Greek						
Active Voice / Future Tenses						
	Εξακολουθητικός Μέλλοντας ( Future Continuous)		Στιγμιαίος Μέλλοντας (Simple Future)		Συντελεσμένος Μέλλοντας (Future Perfect/Imperfect)	
<b>Singular</b>	θα διαβάζομαι	(I) will be being reading	θα διαβαστώ	(I) will be read	θα έχω διαβαστεί	(I) will have been read
	θα διαβάζεσαι	(you) will be being read	θα διαβαστείς	(you) will be read	θα έχεις διαβαστεί	(you) will have been read
	θα διαβάζεται	(he) will be being read	θα διαβαστεί	(he) will be read	θα έχει διαβαστεί	(he) will have been read
<b>Plural</b>	θα διαβαζόμαστε	(we) will be being read	θα διαβαστούμε	(we) will be read	θα έχουμε διαβαστεί	(we) will have been read
	θα διαβαζόσαστε	(you) will be being read	θα διαβαστείτε	(you) will be read	θα έχετε διαβαστεί	(you) will have been read
	θα διαβάζονται	(they) will be being read	θα διαβαστούν	(they) will be read	θα έχουν διαβαστεί	(they) will have been read

## VII

Verbs of 2<sup>nd</sup> Conjugation Class

The verb “to love” in Greek						
Active Voice / Past Tenses						
	Υπερσυντέλικος (Past Perfect)		Αόριστος (Simple Past)		Παρατατικός (Past continuous )	
<b>Singular</b>	είχα αγαπήσει	(I) had loved	αγάπησα	(I) loved	αγαπούσα (αγάπαγα)	(I) was loving
	είχες αγαπήσει	(you) had loved	αγάπησες	(you) loved	αγαπούσες (αγάπαγες)	(you) were loving
	είχε αγαπήσει	(he) had loved	αγάπησε	(he) loved	αγαπούσε (αγάπαγε)	(he) was loving
<b>Plural</b>	είχαμε αγαπήσει	(we) had loved	αγαπήσαμε	(we) loved	αγαπούσαμε (αγαπάγαμε)	(we) were loving
	είχατε αγαπήσει	(you) had loved	αγαπήσατε	(you) loved	αγαπούσατε (αγαπάγατε)	(you) were loving
	είχαν αγαπήσει	(they) had loved	αγάπησαν	(they) loved	αγαπούσαν (αγάπαγαν)	(they) were loving

## VIII

The verb "to love" in Greek				
Active Voice / Present Tenses				
	Ενεστώτας (Present)		Παρακείμενος (Present Perfect)	
<b>Singular</b>	αγαπώ (αγαπάω)	(I) love /am loving	έχω αγαπήσει	(I) have loved /have been loving
	αγαπάς	(you) love / are loving	έχει αγαπήσει	(you) have loved/have been loving
	αγαπά (αγαπάει)	(he) loves/ is loving	έχει αγαπήσει	(he) has loved/has been loving
<b>Plural</b>	αγαπούμε (αγαπάμε)	(we) love / are loving	έχουμε αγαπήσει	(we) have loved/have been loving
	αγαπάτε	(you) love / are loving	έχετε αγαπήσει	(you) have loved/have been loving
	αγαπάνε	(they) love /are loving	έχουν αγαπήσει	(they) have loved/have been loving

## IX

The verb “to love” in Greek						
Active Voice / Future Tenses						
	Εξακολουθητικός Μέλλοντας ( Future Continuous)		Στιγμιαίος Μέλλοντας (Simple Future)		Συντελεσμένος Μέλλοντας (Future Perfect/Imperfect)	
<b>Singular</b>	θα αγαπώ (θα αγαπάω)	(I) will be loving	θα αγαπήσω	(I) will love	θα έχω αγαπήσει	(I) will have loved /have been loving
	θα αγαπάς	(you) will be loving	θα αγαπήσεις	(you) will love	θα έχει αγαπήσει	(you) will have loved/have been loving
	θα αγαπάει	(he) will be loving	θα αγαπήσει	(he) will love	θα έχει αγαπήσει	(he) will has loved/has been loving
<b>Plural</b>	θα αγαπούμε (θα αγαπάμε)	(we )will loving	θα αγαπήσουμε	(we) will love	θα έχουμε αγαπήσει	(we) will have loved/have been loving
	θα αγαπάτε	(you)will be loving	θα αγαπήσετε	(you) will love	θα έχετε αγαπήσει	(you) will have loved/have been loving
	θα αγαπούν (θα αγαπάνε)	(they) will be loving	θα αγαπήσουν	(they) will love	θα έχουν αγαπήσει	(they) will have loved/have been loving

The verb “to love” in Greek						
Passive Voice / Past Tenses						
	Υπερσυντέλικος (Past Perfect)		Αόριστος (Simple Past)		Παρατατικός (Past continuous )	
<b>Singular</b>	είχα αγαπηθεί	(I) had been loved	αγαπήθηκα	(I) was loved	αγαπιόμουν	(I) was being loved
	είχες αγαπηθεί	(you) had been loved	αγαπήθηκες	(you) were loved	αγαπιόσουν	(you) were being loved
	είχε αγαπηθεί	(he) had been loved	αγαπήθηκε	(he) was loved	αγαπιόταν	(he) was being loved
<b>Plural</b>	είχαμε αγαπηθεί	(we) had been loved	αγαπηθήκαμε	(we) were loved	αγαπιόμαστε (αγαπιόμασταν)	(we) were being loved
	είχατε αγαπηθεί	(you) had been loved	αγαπηθήκατε	(you) were loved	αγαπιόσαστε (αγαπιόσασταν)	(you) were being loved
	είχαν αγαπηθεί	(they) had been loved	αγαπήθηκαν	(they) were loved	αγαπιόνταν (αγαπιόντουσαν)	(they) were being loved

The verb "to love" in Greek				
Passive Voice / Present Tenses				
	Ενεστώτας (Present)		Παρακείμενος (Present Perfect)	
<b>Singular</b>	αγαπιέμαι	(I) am being loved	έχω αγαπηθεί	(I) have been loved
	αγαπιέσαι	(you) are being loved	έχεις αγαπηθεί	(you) have been loved
	αγαπιέται	(he) is being loved	έχει αγαπηθεί	(he) has been loved
<b>Plural</b>	αγαπιόμαστε	(we) are being loved	έχουμε αγαπηθεί	(we) have been loved
	αγαπιόσαστε	(you) are being loved	έχετε αγαπηθεί	(you) have been loved
	αγαπιούνται	(they) are being loved	έχουν αγαπηθεί	(they) have been loved

XII

The verb "to love" in Greek						
Passive Voice / Future Tenses						
	Εξακολουθητικός Μέλλοντας ( Future Continuous)		Στιγμιαίος Μέλλοντας (Simple Future)		Συντελεσμένος Μέλλοντας (Future Perfect/Imperfect)	
<b>Singular</b>	θα αγαπιέμαι	(I) will be loved	θα αγαπηθώ	(I) will love	θα έχω αγαπηθεί	(I) will have been loved
	θα αγαπιέσαι	(you) will be loved	θα αγαπηθείς	(you) will love	θα έχεις αγαπηθεί	(you) will have been loved
	θα αγαπιέται	(he) will be loved	θα αγαπηθεί	(he) will love	θα έχει αγαπηθεί	(he) will have been loved
<b>Plural</b>	θα αγαπιόμαστε	(we) will be loved	θα αγαπηθούμε	(we) will love	θα έχουμε αγαπηθεί	(we) will have been loved
	θα αγαπιόσαστε	(you) will be loved	θα αγαπηθείτε	(you) will love	θα έχετε αγαπηθεί	(you) will have been loved
	θα αγαπιούνται	(they) will be loved	θα αγαπηθούν	(they) will love	θα έχουν αγαπηθεί	(they) will have been loved



**Appendix B: Evaluation of Modified Algorithm Testing Output**

Number of words: 562242

Number of Stems: 138629

Frake's statistics:

Mean number of words per conflation class : 4.05573148475

Index compression factor : 0.75343535346

Unchanged Words : 7309 -----> ratio 0.0130720262642 (unchanged words /total words)

Mean Modified Humming Distance 2.44133602915

Median Modified Humming Distance: 2

<b>Evaluation Output (part)</b>			
The number next to every word notes the modified Hamming distance according to Frakes from its stem			
-----ΑΒΑΘ----- is a stem for 12 words	-----ΚΟΠΑΝΙΖ----- is a stem for 20 words	-----ΠΡΟΠΩΛ----- is a stem for 31 words	-----ΠΟΥΛ----- is a stem for 50 words
ΑΒΑΘΩΝ 2	ΚΟΠΑΝΙΖΩ 1	ΠΡΟΠΩΛΩΝΤΑΣ 5	ΠΟΥΛΩΝΤΑΣ 5
ΑΒΑΘΟΥΣ 3	ΚΟΠΑΝΙΖΟΥΝ 3	ΠΡΟΠΩΛΩ 1	ΠΟΥΛΩ 1
ΑΒΑΘΟΥ 21	ΚΟΠΑΝΙΖΟΥΜΕ 4	ΠΡΟΠΩΛΟΥΣΕΣ 5	ΠΟΥΛΟΥΣΕΣ 5
ΑΒΑΘΟΣ 2	ΚΟΠΑΝΙΖΟΣΟΥΝ 5	ΠΡΟΠΩΛΟΥΣΕ 4	ΠΟΥΛΟΥΣΕ 4
ΑΒΑΘΟΙ 2	ΚΟΠΑΝΙΖΟΣΑΣΤΕ 6	ΠΡΟΠΩΛΟΥΣΑΤΕ 6	ΠΟΥΛΟΥΣΑΤΕ 6
ΑΒΑΘΟ 1	ΚΟΠΑΝΙΖΟΝΤΑΣ 5	ΠΡΟΠΩΛΟΥΣΑΜΕ 6	ΠΟΥΛΟΥΣΑΜΕ 6
ΑΒΑΘΗΣ 2	ΚΟΠΑΝΙΖΟΝΤΑΙ 5	ΠΡΟΠΩΛΟΥΣΑ 4	ΠΟΥΛΟΥΣΑ 4
ΑΒΑΘΗ 1	ΚΟΠΑΝΙΖΟΜΑΣΤΕ 6	ΠΡΟΠΩΛΟΥΝΤΑΙ 6	ΠΟΥΛΟΥΝΤΑΙ 6
ΑΒΑΘΗΣ 2	ΚΟΠΑΝΙΖΟΜΑΙ 4	ΠΡΟΠΩΛΟΥΝ 3	ΠΟΥΛΟΥΝ 3
ΑΒΑΘΗ 1	ΚΟΠΑΝΙΖΕΤΕ 3	ΠΡΟΠΩΛΟΥΜΕ 4	ΠΟΥΛΟΥΜΕ 4
ΑΒΑΘΕΣ 2	ΚΟΠΑΝΙΖΕΤΑΙ 4	ΠΡΟΠΩΛΟΥΜΑΣΤΕ 7	ΠΟΥΛΙΩΝ 3
ΑΒΑΘΕΙΣ 3	ΚΟΠΑΝΙΖΕΣΤΕ 4	ΠΡΟΠΩΛΟΥΜΑΙ 5	ΠΟΥΛΙΟΥΝΤΑΙ 7
ΑΒΑΘΕ 1	ΚΟΠΑΝΙΖΕΣΑΙ 4	ΠΡΟΠΩΛΗΣΩ 3	ΠΟΥΛΙΟΥ 3
ΑΒΑΘΑ 1	ΚΟΠΑΝΙΖΕΣ 2	ΠΡΟΠΩΛΗΣΤΕ 4	ΠΟΥΛΙΟΝΤΑΙ 6
	ΚΟΠΑΝΙΖΕΙΣ 3	ΠΡΟΠΩΛΗΣΟΥ 4	ΠΟΥΛΙΕΤΑΙ 5
	ΚΟΠΑΝΙΖΕΙ 2	ΠΡΟΠΩΛΗΣΕΤΕ 5	ΠΟΥΛΙΕΣΤΕ 5
	ΚΟΠΑΝΙΖΕ 1	ΠΡΟΠΩΛΗΣΕ 3	ΠΟΥΛΙΕΣΑΙ 5
-----ΑΒΑΘΟΥΛΩΤ----- is a stem for 11 words	ΚΟΠΑΝΙΖΑΝ 2	ΠΡΟΠΩΛΗΣΑΤΕ 5	ΠΟΥΛΙΕΣ 3
	ΚΟΠΑΝΙΖΑΜΕ 3	ΠΡΟΠΩΛΗΣΑΜΕ 5	ΠΟΥΛΙΕΜΑΙ 5

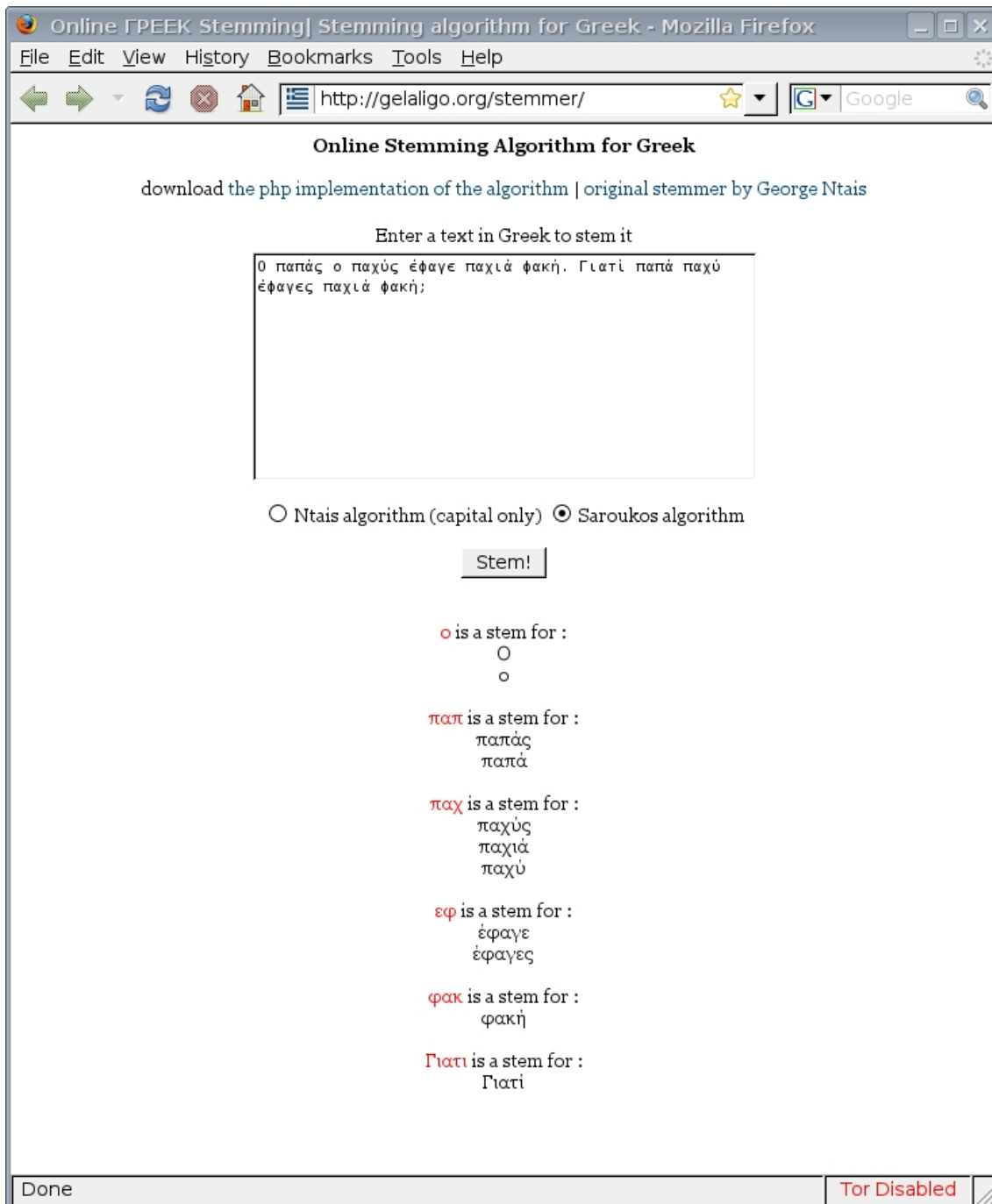
XIV

ΑΒΑΘΟΥΛΩΤΩΝ 2	ΚΟΠΑΝΙΖΑ 1	ΠΡΟΠΩΛΗΣΑ 3	ΠΟΥΛΙΑΣ 3
ΑΒΑΘΟΥΛΩΤΟΥΣ 2	-----ΚΟΠΑΝΙΖΟΤ----- is a stem for 1 words	ΠΡΟΠΩΛΗΜΕΝΟΣ 6	ΠΟΥΛΙΑ 2
ΑΒΑΘΟΥΛΩΤΟΥ 2		ΠΡΟΠΩΛΗΜΕΝΟ 5	<b>ΠΟΥΛΙ</b>
ΑΒΑΘΟΥΛΩΤΟΣ 2	ΚΟΠΑΝΙΖΟΤΑΝ 2	ΠΡΟΠΩΛΗΜΕΝΗ 5	ΠΟΥΛΗΣΩ 3
ΑΒΑΘΟΥΛΩΤΟΙ 2		ΠΡΟΠΩΛΗΜΕΝΕ 5	ΠΟΥΛΗΣΤΕ 4
ΑΒΑΘΟΥΛΩΤΟ 1	-----ΚΟΠΑΝΙΖΟΣΑΣΤ----- is a stem for 1 words	ΠΡΟΠΩΛΗΘΩ 3	ΠΟΥΛΗΣΟΥ 4
ΑΒΑΘΟΥΛΩΤΗΣ 2		ΠΡΟΠΩΛΗΘΗΚΕΣ 6	ΠΟΥΛΗΣΕΤΕ 5
ΑΒΑΘΟΥΛΩΤΗ 1	ΚΟΠΑΝΙΖΟΣΑΣΤΑΝ 2	ΠΡΟΠΩΛΗΘΗΚΕ 5	ΠΟΥΛΗΣΕ 3
ΑΒΑΘΟΥΛΩΤΕΣ 2		ΠΡΟΠΩΛΗΘΗΚΑ 5	ΠΟΥΛΗΣΑΤΕ 5
ΑΒΑΘΟΥΛΩΤΕ 1	-----ΚΟΠΑΝΙΖΟΝΤΟΥΣ----- is a stem for 1 words	ΠΡΟΠΩΛΕΙΤΕ 4	ΠΟΥΛΗΣΑΜΕ 5
ΑΒΑΘΟΥΛΩΤΑ 1		ΠΡΟΠΩΛΕΙΣ 3	ΠΟΥΛΗΣΑ 3
	ΚΟΠΑΝΙΖΟΝΤΟΥΣΑΝ 2	ΠΡΟΠΩΛΕΙ 2	ΠΟΥΛΗΜΕΝΟΣ 6
-----ΑΒΑΘΜΟΛΟΓΗΤ----- is a stem for 11 words		-----ΠΡΟΠΩΛΟΥΤ----- is a stem for 1 words	ΠΟΥΛΗΜΕΝΟ 5
	-----ΚΟΠΑΝΙΖΟΝΤ----- is a stem for 1 words	ΠΡΟΠΩΛΟΥΤΑΝ 2	ΠΟΥΛΗΜΕΝΗ 5
ΑΒΑΘΜΟΛΟΓΗΤΩΝ 2	ΚΟΠΑΝΙΖΟΝΤΑΝ 2	-----ΠΡΟΠΩΛΟΥΣ----- is a stem for 2 words	ΠΟΥΛΗΜΕΝΕ 5
ΑΒΑΘΜΟΛΟΓΗΤΟΥΣ 3		ΠΡΟΠΩΛΟΥΣΟΥΝ 3	ΠΟΥΛΗΘΩ 3
ΑΒΑΘΜΟΛΟΓΗΤΟΥ 2	-----ΚΟΠΑΝΙΖΟΜ----- is a stem for 1 words	ΠΡΟΠΩΛΟΥΣΑΝ 2	ΠΟΥΛΗΘΗΚΕΣ 6
ΑΒΑΘΜΟΛΟΓΗΤΟΣ 2		-----	ΠΟΥΛΗΘΗΚΕ 5
ΑΒΑΘΜΟΛΟΓΗΤΟΙ 2	ΚΟΠΑΝΙΖΟΜΟΥΝ 3	ΠΡΟΠΩΛΟΥΣΑΣΤ----- is a stem for 1 words	ΠΟΥΛΗΘΗΚΑ 5
ΑΒΑΘΜΟΛΟΓΗΤΟ 1	-----ΚΟΠΑΝΙΖΟΜΑΣΤ----- is a stem for 1 words	ΠΡΟΠΩΛΟΥΣΑΣΤΑΝ 2	ΠΟΥΛΑΩ 2
ΑΒΑΘΜΟΛΟΓΗΤΗΣ 2		-----ΠΡΟΠΩΛΟΥΝΤ----- is a stem for 1 words	ΠΟΥΛΑΣ 2
ΑΒΑΘΜΟΛΟΓΗΤΗ 1	ΚΟΠΑΝΙΖΟΜΑΣΤΑΝ 2		ΠΟΥΛΑΝΕ 3
ΑΒΑΘΜΟΛΟΓΗΤΕΣ 2	-----ΚΟΠΑΝΙΖΑΤ----- is a stem for 1 words	ΠΡΟΠΩΛΟΥΝΤΑΝ 2	ΠΟΥΛΑΝ 2
ΑΒΑΘΜΟΛΟΓΗΤΕ 1		-----	ΠΟΥΛΑΜΕ 3
ΑΒΑΘΜΟΛΟΓΗΤΑ 1	ΚΟΠΑΝΙΖΑΤΕ 1	ΠΡΟΠΩΛΟΥΜΑΣΤ----- is a stem for 1 words	ΠΟΥΛΑΕΙ 3
	-----ΚΟΠΑΝΗΣ----- is a stem for 6 words	ΠΡΟΠΩΛΟΥΜΑΣΤΑΝ 2	ΠΟΥΛΑΓΕΣ 4
-----ΑΒΑΘΜΙΔΩΤ----- is a stem for 11 words		-----	ΠΟΥΛΑΓΕ 3
	ΚΟΠΑΝΗΣΟΥΝ 3	ΠΡΟΠΩΛΟΥΜΑΣΤΑΝ 2	ΠΟΥΛΑΓΑΤΕ 5
ΑΒΑΘΜΙΔΩΤΩΝ 2	ΚΟΠΑΝΗΣΟΥΜΕ 4	-----ΠΡΟΠΩΛΗΣ----- is a stem for 10 words	ΠΟΥΛΑΓΑΝΕ 5
ΑΒΑΘΜΙΔΩΤΟΥΣ 3	ΚΟΠΑΝΗΣΕΣ 2		ΠΟΥΛΑΓΑΜΕ 5
ΑΒΑΘΜΙΔΩΤΟΥ 2	ΚΟΠΑΝΗΣΕΙΣ 3	ΠΡΟΠΩΛΗΣΟΥΝ 3	ΠΟΥΛΑΓΑ 3
ΑΒΑΘΜΙΔΩΤΟΣ 2	ΚΟΠΑΝΗΣΕΙ 2	ΠΡΟΠΩΛΗΣΟΥΜΕ 4	ΠΟΥΛΑ 1
ΑΒΑΘΜΙΔΩΤΟΙ 2	ΚΟΠΑΝΗΣΑΝ 2	ΠΡΟΠΩΛΗΣΗΣ 2	-----
ΑΒΑΘΜΙΔΩΤΟ 1		ΠΡΟΠΩΛΗΣΗ 1	ΠΟΥΛΑΧΕΡΙΑΡ----- is a stem for 1 words
ΑΒΑΘΜΙΔΩΤΗΣ 2		ΠΡΟΠΩΛΗΣΕΩΣ 3	ΠΟΥΛΑΧΕΡΙΑΡ 0
ΑΒΑΘΜΙΔΩΤΗ 1		ΠΡΟΠΩΛΗΣΕΩΝ 3	-----
ΑΒΑΘΜΙΔΩΤΕΣ 2			ΠΟΥΛΑΧΕΡ----- is a stem for 1 words
ΑΒΑΘΜΙΔΩΤΕ 1			ΠΟΥΛΑΧΕΡΙΑ 2
ΑΒΑΘΜΙΔΩΤΑ 1			-----
			ΠΟΥΛΟΥΣ----- is a stem for 1 words
			ΠΟΥΛΟΥΣΑΝ 2

<p>-----ΑΒΑΕΙ----- is a stem for 1 words</p> <p>ΑΒΑΕΙΟ 1</p>		<p>ΠΡΟΠΩΛΗΣΕΣ 2 ΠΡΟΠΩΛΗΣΕΙΣ 3 ΠΡΟΠΩΛΗΣΕΙ 2 ΠΡΟΠΩΛΗΣΑΝ 2</p>	<p>----- ΠΟΥΛΟΛΟΓ----- is a stem for 1 words</p> <p>ΠΟΥΛΟΛΟΓΟΣ 2</p> <p>----- ΠΟΥΛΟΒΕΡΑΚ----- - is a stem for 2 words</p> <p>ΠΟΥΛΟΒΕΡΑΚΙΑ 2 ΠΟΥΛΟΒΕΡΑΚΙ 1</p> <p>----- ΠΟΥΛΟΒΕΡ----- is a stem for 1 words</p> <p>ΠΟΥΛΟΒΕΡ 0</p>
--	--	---	--

The previous table is a small sample from an extensive list of conflation classes created by the modified stemmer. In this small list we included examples of both understemming and overstemming. One overstemming example is the conflation class of “ΠΟΥΛ” which includes words from the word “πουλί” (bird) and “πουλάω” (to sell). On the contrary, understemming examples have to do with the grammatical changes that words undergo in past tenses.

## Appendix B: On-line Stemmer



Our implementation of the algorithm is freely available at <http://gelaligo.org/stemmer> under an LGPL licence, along with an on-line demo.

**Appendix C: Stop Word List**

ΑΒΑ	ΑΝΑΜΕΣΑ	ΑΥΡΙΟ
ΑΓΑ	ΑΝΑΜΕΤΑΞΥ	ΑΦΗ
ΑΓΗ	ΑΝΕΥ	ΑΦΟΤΟΥ
ΑΓΩ	ΑΝΤΙ	ΑΦΟΥ
ΑΔΗ	ΑΝΤΙΠΕΡΑ	ΑΧ
ΑΔΩ	ΑΝΤΙΣ	ΑΧΕ
ΑΕ	ΑΝΩ	ΑΧΟ
ΑΕΙ	ΑΝΩΤΕΡΩ	ΑΨΑ
ΑΘΩ	ΑΞΑΦΝΑ	ΑΨΕ
ΑΙ	ΑΠ	ΑΨΗ
ΑΙΚ	ΑΠΕΝΑΝΤΙ	ΑΨΥ
ΑΚΗ	ΑΠΟ	ΑΩΕ
ΑΚΟΜΑ	ΑΠΟΨΕ	ΑΩΟ
ΑΚΟΜΗ	ΑΠΩ	ΒΑΝ
ΑΚΡΙΒΩΣ	ΑΡΑ	ΒΑΤ
ΑΛΑ	ΑΡΑΓΕ	ΒΑΧ
ΑΛΗΘΕΙΑ	ΑΡΕ	ΒΕΑ
ΑΛΗΘΙΝΑ	ΑΡΚ	ΒΕΒΑΙΟΤΑΤΑ
ΑΛΛΑΧΟΥ	ΑΡΚΕΤΑ	ΒΗΞ
ΑΛΛΙΩΣ	ΑΡΛ	ΒΙΑ
ΑΛΛΙΩΤΙΚΑ	ΑΡΜ	ΒΙΕ
ΑΛΛΟΙΩΣ	ΑΡΤ	ΒΙΗ
ΑΛΛΟΙΩΤΙΚΑ	ΑΡΥ	ΒΙΟ
ΑΛΛΟΤΕ	ΑΡΩ	ΒΟΗ
ΑΛΤ	ΑΣ	ΒΟΩ
ΑΛΩ	ΑΣΑ	ΒΡΕ
ΑΜΑ	ΑΣΟ	ΓΑ
ΑΜΕ	ΑΤΑ	ΓΑΒ
ΑΜΕΣΑ	ΑΤΕ	ΓΑΡ
ΑΜΕΣΩΣ	ΑΤΗ	ΓΕΝ
ΑΜΩ	ΑΤΙ	ΓΕΣ
ΑΝ	ΑΤΜ	ΓΗ
ΑΝΑ	ΑΤΟ	ΓΗΝ

XVIII

ΓΙ	ΔΥΕ	ΕΞ
ΓΙΑ	ΔΥΟ	ΕΞΑΦΝΑ
ΓΙΕ	ΔΩ	ΕΞΙ
ΓΙΝ	ΕΑΜ	ΕΞΙΣΟΥ
ΓΙΟ	ΕΑΝ	ΕΞΩ
ΓΚΙ	ΕΑΡ	ΕΟΚ
ΓΚΥ	ΕΘΗ	ΕΠΑΝΩ
ΓΟΗ	ΕΙ	ΕΠΕΙΔΗ
ΓΟΟ	ΕΙΔΕΜΗ	ΕΠΕΙΤΑ
ΓΡΗΓΟΡΑ	ΕΙΘΕ	ΕΠΗ
ΓΡΙ	ΕΙΜΑΙ	ΕΠΙ
ΓΡΥ	ΕΙΜΑΣΤΕ	ΕΠΙΣΗΣ
ΓΥΗ	ΕΙΝΑΙ	ΕΠΟΜΕΝΩΣ
ΓΥΡΩ	ΕΙΣ	ΕΡΑ
ΔΑ	ΕΙΣΑΙ	ΕΣ
ΔΕ	ΕΙΣΑΣΤΕ	ΕΣΑΣ
ΔΕΗ	ΕΙΣΤΕ	ΕΣΕ
ΔΕΙ	ΕΙΤΕ	ΕΣΕΙΣ
ΔΕΝ	ΕΙΧΑ	ΕΣΕΝΑ
ΔΕΣ	ΕΙΧΑΜΕ	ΕΣΗ
ΔΗ	ΕΙΧΑΝ	ΕΣΤΩ
ΔΗΘΕΝ	ΕΙΧΑΤΕ	ΕΣΥ
ΔΗΛΑΔΗ	ΕΙΧΕ	ΕΣΩ
ΔΗΩ	ΕΙΧΕΣ	ΕΤΙ
ΔΙ	ΕΚ	ΕΤΣΙ
ΔΙΑ	ΕΚΟ	ΕΥ
ΔΙΑΡΚΩΣ	ΕΚΕΙ	ΕΥΑ
ΔΙΟΛΟΥ	ΕΛΑ	ΕΥΓΕ
ΔΙΣ	ΕΛΙ	ΕΥΘΥΣ
ΔΙΧΩΣ	ΕΜΠ	ΕΥΤΥΧΩΣ
ΔΟΛ	ΕΝ	ΕΦΕ
ΔΟΝ	ΕΝΤΕΛΩΣ	ΕΦΕΞΗΣ
ΔΡΑ	ΕΝΤΟΣ	ΕΦΤ
ΔΡΥ	ΕΝΤΩΜΕΤΑΞΥ	ΕΧΕ
ΔΡΧ	ΕΝΩ	ΕΧΕΙ

XIX

ΕΧΕΙΣ	ΙΒΟ	ΚΑΠΩΣ
ΕΧΕΤΕ	ΙΔΗ	ΚΑΤ
ΕΧΘΕΣ	ΙΔΙΩΣ	ΚΑΤΑ
ΕΧΟΜΕ	ΙΕ	ΚΑΤΙ
ΕΧΟΥΜΕ	ΙΙ	ΚΑΤΙΤΙ
ΕΧΟΥΝ	ΙΙΙ	ΚΑΤΟΠΙΝ
ΕΧΤΕΣ	ΙΚΑ	ΚΑΤΩ
ΕΧΩ	ΙΛΟ	ΚΑΩ
ΕΩΣ	ΙΜΑ	ΚΒΟ
ΖΕΑ	ΙΝΑ	ΚΕΑ
ΖΕΗ	ΙΝΩ	ΚΕΙ
ΖΕΙ	ΙΞΕ	ΚΕΝ
ΖΕΝ	ΙΞΟ	ΚΙ
ΖΗΝ	ΙΟ	ΚΙΜ
ΖΩ	ΙΟΙ	ΚΙΟΛΑΣ
Η	ΙΣΑ	ΚΙΤ
ΗΔΗ	ΙΣΑΜΕ	ΚΙΧ
ΗΔΥ	ΙΣΕ	ΚΚΕ
ΗΘΗ	ΙΣΗ	ΚΛΙΣΕ
ΗΛΟ	ΙΣΙΑ	ΚΛΠ
ΗΜΙ	ΙΣΟ	ΚΟΚ
ΗΠΑ	ΙΣΩΣ	ΚΟΝΤΑ
ΗΣΑΣΤΕ	ΙΩΒ	ΚΟΧ
ΗΣΟΥΝ	ΙΩΝ	ΚΤΑ
ΗΤΑ	ΙΩΣ	ΚΥΡ
ΗΤΑΝ	Ιαυ	ΚΥΡΙΩΣ
ΗΤΑΝΕ	ΚΑΘ	ΚΩ
ΗΤΟΙ	ΚΑΘΕ	ΚΩΝ
ΗΤΤΟΝ	ΚΑΘΕΤΙ	ΛΑ
ΗΩ	ΚΑΘΟΛΟΥ	ΛΕΑ
ΘΑ	ΚΑΘΩΣ	ΛΕΝ
ΘΥΕ	ΚΑΙ	ΛΕΟ
ΘΩΡ	ΚΑΝ	ΛΙΑ
Ι	ΚΑΠΟΤΕ	ΛΙΓΑΚΙ
ΙΑ	ΚΑΠΟΥ	ΛΙΓΟ

ΛΙΓΩΤΕΡΟ	ΜΗΝ	ΝΤΑ
ΛΙΟ	ΜΗΠΩΣ	ΝΤΕ
ΛΙΡ	ΜΗΤΕ	ΝΤΙ
ΛΟΓΩ	ΜΙ	ΝΤΟ
ΛΟΙΠΑ	ΜΙΕ	ΝΥΝ
ΛΟΙΠΟΝ	ΜΙΣ	ΝΩΕ
ΛΟΣ	ΜΜΕ	ΝΩΡΙΣ
ΛΣ	ΜΝΑ	ΞΑΝΑ
ΛΥΩ	ΜΟΒ	ΞΑΦΝΙΚΑ
ΜΑ	ΜΟΛΙΣ	ΞΕΩ
ΜΑΖΙ	ΜΟΛΟΝΟΤΙ	ΞΙ
ΜΑΚΑΡΙ	ΜΟΝΑΧΑ	Ο
ΜΑΛΙΣΤΑ	ΜΟΝΟΜΙΑΣ	ΟΑ
ΜΑΛΛΟΝ	ΜΟΥ	ΟΑΠ
ΜΑΝ	ΜΠΑ	ΟΔΟ
ΜΑΞ	ΜΠΟΡΕΙ	ΟΕ
ΜΑΣ	ΜΠΟΡΟΥΝ	ΟΖΟ
ΜΑΤ	ΜΠΡΑΒΟ	ΟΗΕ
ΜΕ	ΜΠΡΟΣ	ΟΙ
ΜΕΘΑΥΡΙΟ	ΜΠΩ	ΟΙΑ
ΜΕΙ	ΜΥ	ΟΙΗ
ΜΕΙΟΝ	ΜΥΑ	ΟΚΑ
ΜΕΛ	ΜΥΝ	ΟΛΟΓΥΡΑ
ΜΕΛΕΙ	ΝΑ	ΟΛΟΝΕΝ
ΜΕΛΛΕΤΑΙ	ΝΑΕ	ΟΛΟΤΕΛΑ
ΜΕΜΙΑΣ	ΝΑΙ	ΟΛΩΣΔΙΟΛΟΥ
ΜΕΝ	ΝΑΟ	ΟΜΩΣ
ΜΕΣ	ΝΔ	ΟΝ
ΜΕΣΑ	ΝΕῖ	ΟΝΕ
ΜΕΤ	ΝΙ	ΟΝΟ
ΜΕΤΑ	ΝΙΑ	ΟΠΑ
ΜΕΤΑΞΥ	ΝΙΚ	ΟΠΕ
ΜΕΧΡΙ	ΝΙΛ	ΟΠΗ
ΜΗ	ΝΙΝ	ΟΠΟ
ΜΗΔΕ	ΝΙΟ	ΟΠΟΙΑΔΗΠΟΤΕ



ΟΠΟΙΑΝΔΗΠΟΤΕ	ΟΤΕ	ΠΙΚ
ΟΠΟΙΑΣΔΗΠΟΤΕ	ΟΤΙ	ΠΙΟ
ΟΠΟΙΔΗΠΟΤΕ	ΟΤΙΔΗΠΟΤΕ	ΠΙΣΩ
ΟΠΟΙΕΣΔΗΠΟΤΕ	ΟΥ	ΠΙΤ
ΟΠΟΙΟΔΗΠΟΤΕ	ΟΥΔΕ	ΠΙΩ
ΟΠΟΙΟΝΔΗΠΟΤΕ	ΟΥΚ	ΠΛΑΙ
ΟΠΟΙΟΣΔΗΠΟΤΕ	ΟΥΣ	ΠΛΕΟΝ
ΟΠΟΙΟΥΔΗΠΟΤΕ	ΟΥΤΕ	ΠΛΗΝ
ΟΠΟΙΟΥΣΔΗΠΟΤΕ	ΟΥΦ	ΠΛΩ
ΟΠΟΙΩΝΔΗΠΟΤΕ	ΟΧΙ	ΠΜ
ΟΠΟΤΕΔΗΠΟΤΕ	ΟΨΑ	ΠΟΑ
ΟΠΟΥ	ΟΨΕ	ΠΟΕ
ΟΠΟΥΔΗΠΟΤΕ	ΟΨΗ	ΠΟΛ
ΟΠΩΣ	ΟΨΙ	ΠΟΛΥ
ΟΡΑ	ΟΨΟ	ΠΟΠ
ΟΡΕ	ΠΑ	ΠΟΤΕ
ΟΡΗ	ΠΑΛΙ	ΠΟΥ
ΟΡΟ	ΠΑΝ	ΠΟΥΘΕ
ΟΡΦ	ΠΑΝΤΟΤΕ	ΠΟΥΘΕΝΑ
ΟΡΩ	ΠΑΝΤΟΥ	ΠΡΕΠΕΙ
ΟΣΑ	ΠΑΝΤΩΣ	ΠΡΙ
ΟΣΑΔΗΠΟΤΕ	ΠΑΠ	ΠΡΙΝ
ΟΣΕ	ΠΑΡ	ΠΡΟ
ΟΣΕΣΔΗΠΟΤΕ	ΠΑΡΑ	ΠΡΟΚΕΙΜΕΝΟΥ
ΟΣΗΔΗΠΟΤΕ	ΠΕΙ	ΠΡΟΚΕΙΤΑΙ
ΟΣΗΝΔΗΠΟΤΕ	ΠΕΡ	ΠΡΟΠΕΡΣΙ
ΟΣΗΣΔΗΠΟΤΕ	ΠΕΡΑ	ΠΡΟΣ
ΟΣΟΔΗΠΟΤΕ	ΠΕΡΙ	ΠΡΟΤΟΥ
ΟΣΟΙΔΗΠΟΤΕ	ΠΕΡΙΠΟΥ	ΠΡΟΧΘΕΣ
ΟΣΟΝΔΗΠΟΤΕ	ΠΕΡΣΙ	ΠΡΟΧΤΕΣ
ΟΣΟΣΔΗΠΟΤΕ	ΠΕΡΥΣΙ	ΠΡΩΤΥΤΕΡΑ
ΟΣΟΥΔΗΠΟΤΕ	ΠΕΣ	ΠΥΑ
ΟΣΟΥΣΔΗΠΟΤΕ	ΠΙ	ΠΥΞ
ΟΣΩΝΔΗΠΟΤΕ	ΠΙΑ	ΠΥΟ
ΟΤΑΝ	ΠΙΘΑΝΟΝ	ΠΥΡ

## XXII

ΠΧ	ΣΕΡ	ΤΑ
ΠΩ	ΣΕΤ	ΤΑΔΕ
ΠΩΛ	ΣΕΦ	ΤΑΚ
ΠΩΣ	ΣΗΜΕΡΑ	ΤΑΝ
ΡΑ	ΣΙ	ΤΑΟ
ΡΑΙ	ΣΙΑ	ΤΑΥ
ΡΑΠ	ΣΙΓΑ	ΤΑΧΑ
ΡΑΣ	ΣΙΚ	ΤΑΧΑΤΕ
ΡΕ	ΣΙΧ	ΤΕ
ΡΕΑ	ΣΚΙ	ΤΕΙ
ΡΕΕ	ΣΟΙ	ΤΕΛ
ΡΕΙ	ΣΟΚ	ΤΕΛΙΚΑ
ΡΗΣ	ΣΟΛ	ΤΕΛΙΚΩΣ
ΡΘΩ	ΣΟΝ	ΤΕΣ
ΡΙΟ	ΣΟΣ	ΤΕΤ
ΡΟ	ΣΟΥ	ΤΖΟ
ΡΟῖ	ΣΡΙ	ΤΗ
ΡΟΕ	ΣΤΑ	ΤΗΛ
ΡΟΖ	ΣΤΗ	ΤΗΝ
ΡΟΗ	ΣΤΗΝ	ΤΗΣ
ΡΟΘ	ΣΤΗΣ	ΤΙ
ΡΟΙ	ΣΤΙΣ	ΤΙΚ
ΡΟΚ	ΣΤΟ	ΤΙΜ
ΡΟΛ	ΣΤΟΝ	ΤΙΠΟΤΑ
ΡΟΝ	ΣΤΟΥ	ΤΙΠΟΤΕ
ΡΟΣ	ΣΤΟΥΣ	ΤΙΣ
ΡΟΥ	ΣΤΩΝ	ΤΝΤ
ΣΑΙ	ΣΥ	ΤΟ
ΣΑΝ	ΣΥΓΧΡΟΝΩΣ	ΤΟΙ
ΣΑΟ	ΣΥΝ	ΤΟΚ
ΣΑΣ	ΣΥΝΑΜΑ	ΤΟΜ
ΣΕ	ΣΥΝΕΠΩΣ	ΤΟΝ
ΣΕΙΣ	ΣΥΝΗΘΩΣ	ΤΟΠ
ΣΕΚ	ΣΧΕΔΟΝ	ΤΟΣ
ΣΕΞ	ΣΩΣΤΑ	ΤΟΣΩΝ

## XXIII

ΤΟΣΑ	ΥΠΟΨΙΝ	ΧΟΗ
ΤΟΣΕΣ	ΥΣΤΕΡΑ	ΧΟΛ
ΤΟΣΗ	ΥΦΗ	ΧΡΩ
ΤΟΣΗΝ	ΥΨΗ	ΧΤΕΣ
ΤΟΣΗΣ	ΦΑ	ΧΩΡΙΣ
ΤΟΣΟ	ΦΑῖ	ΧΩΡΙΣΤΑ
ΤΟΣΟΙ	ΦΑΕ	ΨΕΣ
ΤΟΣΟΝ	ΦΑΝ	ΨΗΛΑ
ΤΟΣΟΣ	ΦΑΞ	ΨΙ
ΤΟΣΟΥ	ΦΑΣ	ΨΙΤ
ΤΟΣΟΥΣ	ΦΑΩ	Ω
ΤΟΤΕ	ΦΕΖ	ΩΑ
ΤΟΥ	ΦΕΙ	ΩΑΣ
ΤΟΥΛΑΧΙΣΤΟ	ΦΕΤΟΣ	ΩΔΕ
ΤΟΥΛΑΧΙΣΤΟΝ	ΦΕΥ	ΩΕΣ
ΤΟΥΣ	ΦΙ	ΩΘΩ
ΤΣ	ΦΙΑ	ΩΜΑ
ΤΣΑ	ΦΙΣ	ΩΜΕ
ΤΣΕ	ΦΟΞ	ΩΝ
ΤΥΧΟΝ	ΦΠΑ	ΩΟ
ΤΩ	ΦΡΙ	ΩΟΝ
ΤΩΝ	ΧΑ	ΩΟΥ
ΤΩΡΑ	ΧΑΗ	ΩΣ
ΥΑΣ	ΧΑΛ	ΩΣΑΝ
ΥΒΑ	ΧΑΝ	ΩΣΗ
ΥΒΟ	ΧΑΦ	ΩΣΟΤΟΥ
ΥΙΕ	ΧΕ	ΩΣΠΟΥ
ΥΙΟ	ΧΕΙ	ΩΣΤΕ
ΥΛΑ	ΧΘΕΣ	ΩΣΤΟΣΟ
ΥΛΗ	ΧΙ	ΩΤΑ
ΥΝΙ	ΧΙΑ	ΩΧ
ΥΠ	ΧΙΑ	ΩΩΝ
ΥΠΕΡ	ΧΙΟ	ΓΙΑΤΙ
ΥΠΟ	ΧΛΜ	
ΥΠΟΨΗ	ΧΜ	