

Constructing pseudowords for experimental research: Problems and solutions*

Anthi Revithiadou, Dimitra Ioannou, Maria Chatzinikolaou & Katerina

Aivazoglou

Aristotle University of Thessaloniki

revith@lit.auth.gr, dioannos@lit.auth.gr, cnmaria@lit.auth.gr, aaivazog@lit.auth.gr

Abstract

In experimental research, it is often the case that certain theoretical assumptions and hypotheses are tested on the basis of a set of constructed/novel data. In this article, we present the methodology for the construction of pseudowords for the purposes of an experiment that aims at testing the default position of stress in Greek nouns. In order to achieve this goal, we relied on corpora and other associated quantitative tools which are constructed exclusively for Greek by ILSP and are freely available on the web (<http://speech.ilsp.gr/iplr>). Our main goal was to create 200 pseudonouns that would sound native enough but yet not too familiar to the native speakers in a measurable way. For this purpose, we manipulated the available sources by creating a finer-grained corpus which incorporates information on the morphological category, word size and syllable type of its items.

Key words: pseudowords, experimental research, default stress, Clean Corpus

* This research was supported by the ‘Excellence 2011’ Program [Action B, Ref.No 87883] awarded to Dr. Anthi Revithiadou and funded by the Research Committee of the Aristotle University of Thessaloniki.

1. Setting the stage: Stress in Greek

Greek language has morphology-determined stress, that is, stress cannot be accounted for on merely phonological grounds but is assigned on the basis of grammar-specific principles. More specifically, the majority of morphemes, such as stems/roots, thematic vowels, derivational suffixes and inflectional endings have lexically-encoded accentual properties (i.e., they are accented, post-/pre-accenting) (Revithiadou 1999). Given that words in Greek consist of more than one morpheme, it is often the case that morphemes with different accentual properties may coexist in a word. Revithiadou (1999) has argued that when there is a conflict between accented morphemes, morphology offers a helping hand in deciding which accented morpheme will win. In the absence of lexically-encoded stress information, however, stress is on the syllable dictated by the language-specific default, that is, on the antepenultimate, e.g., *píthikos* ‘monkey’ (Malikouti-Drachman & Drachman 1989; Ralli & Touratzidis 1992; Revithiadou 1999, 2007; Burzio & Tantalou 2007, a.o.). Morphology-oriented stress in combination with boundedness to the last three syllables of the word yield only three possible stress patterns for Greek: antepenultimate (APU), penultimate (PU) and ultimate (U) stress (Malikouti-Drachman & Drachman 1989; Drachman & Malikouti-Drachman 1999):

- | | | | | |
|-----|----|-----|----------|---------------------|
| (1) | a. | APU | píthikos | ‘monkey-NOM.SG’ |
| | b. | PU | tsobános | ‘shepherd- NOM.SG’ |
| | c. | U | maragós | ‘carpenter- NOM.SG’ |

- (2) a. APU γίtonas ‘neighbor-NOM.SG’
 b. PU eónas ‘century-NOM.SG’
 c. U vasiljás ‘king-NOM.SG’
- (3) a. APU γέfira ‘bridge-NOM.SG’
 b. PU elpída ‘hope-NOM.SG’
 c. U aγorá ‘market-NOM.SG’

The examples in (4b-c) demonstrate lexically-inflected stress patterns, whereas the example in (4a), in which an accentless morpheme is combined with an accentless inflectional ending, represents the *phonological default* (PDf).

- (4) a. /γiton-as/ accentless root
 b. /eón-as/ accented root
 c. /vasilj[^]-as/ post-accenting root ([^] = non-local accent)

It is important to note that PDf is an analysis-specific construct, which means that a differentiation in the analysis may result to a different definition of the PDf within the same language. For example in Russian, two different patterns have been proposed to represent the language PDf: (a) default is word initial (Halle 1973, 1997; Kiparsky & Halle 1977; Melvold 1990) and (b) default is post-stem (Alderete 1999, 2001a,b).

Although APU has been acknowledged to represent the phonological aspect of the Greek stress system (Malikouti-Drachman & Drachman 1989; Ralli & Touratzidis 1992; Revithiadou 1999, 2007, Burzio & Tantalou 2007 a.o.), it has not been experimentally shown that is also the most prevalent or statistically preferred pattern.

On the contrary, Protopapas et al. (2006) showed that it is the least preferred in reading tasks, whereas a number of recent studies have shown that APU stress is marginal in suffixless words, e.g. acronyms (see Nikolou et al. 2012; Topintzi & Kainada 2012), and in certain classes of inflected words, e.g. nouns in *-a* (Apostolouda 2012). In order to shed light on this issue, we designed a production and a perception experiment that aimed at exploring (a) whether a stress system with a profoundly morphological stress has a robustly present PDF and (b) whether morphological classhood functions as a cue for stress. Our research questions focus exclusively on the stress behavior of nouns.

For the purposes of the experiment, we were required to construct 200 pseudowords from five (5) major morphological classes: nouns ending in *-os*, *-o*, *-a*, *-as*, and *-i* (fem/neut), and of specific size and syllable structure. The pseudowords were constructed on the basis of actual words (5a). The segmental (C, V) positions that were subject to modification are underlined in (5b).

	<i>actual words</i>	<i>target C/V positions</i>
(5) a.	CCV.CV.CV	b. CCV. <u>CV</u> .CV(C)
	CV.CV.CV	CV. <u>CV</u> .CV(C)
	CCV.CV	<u>CCV</u> .CV(C)
	CV.CV	<u>CV</u> .CV(C)

Our main concern for the data sets (i.e., pseudowords) used at the experiment was that they should be constructed in a way that they were not familiar but still ‘sounded’ Greek enough to the Greek speakers’ ears. For this purpose, we developed a specific methodology that exploits corpus-based tools that are freely available on the web, but

at the same time takes into consideration the morphological classhood and the morphosyntactic category of words (i.e., nouns).

3. The methodology of pseudoword construction

3.1. Constructing a category-specific corpus

The construction of pseudowords for a production experiment on morphology-oriented stress was proven a quite challenging task. The major methodological issue was to find a reliable way to measure the degree of the familiarity of the pseudowords or, in simpler words, to make the speakers' 'familiarity intuition' measurable. In order to achieve this goal, we relied on an existing corpus, namely *Clean*. The *Clean Corpus* is created by Protopapas and his colleagues and is a component of the "ILSP Psycho-Linguistic Resource" (<http://speech.ilsp.gr/iplr>, cf. Protopapas et al. 2010). It is a medium size corpus which contains over 200.000 types (29.000.000 tokens)¹ culled up mainly from newspapers, magazines, etc.

The major advantage of *Clean* is that it is freely accessible on the web and, more importantly, provides a set of quantitative measures for each word. The variables that were deemed relevant for our study are the following:

- (6) a. Bigram frequencies (phonemes only): i. Logmean bigram token frequency; ii. Logmean bigram type frequency.
- b. Neighborhoods & cohorts: i. N phonological neighbors (replace only); ii. N phonological neighbors (replace, delete, insert, transpose); iii. Phonological Levenshtein distance 20.

¹ A *type* is the unique form of a word, while a *token* is any occurrence of that particular word.

Bigrams are pairs of adjacent items; in phonological representations bigrams refer to pairs of phonemes:²

Bigram counts are calculated by first summing up all the occurrences (tokens) of each combination of two phones. Total bigram frequency is related to the difficulty with which an item can be read, as it reflects the familiarity of the reader with the combinations of phones exhibited by a given item (word) [...].

The neighbors of an item are items (words) of equal length that differ from the probe item by a single segment.

(URL: <http://speech.ilsp.gr/iplr/documentation.htm>)

There were several variables concerning bigrams. We excluded the variables that were taking into account the orthographical representation of the word and focused on the ones that take into account the phonological representation. More specifically, we chose (i) Logmean bigram token frequency and (ii) Logmean bigram type frequency, which focus only on phonemes of tokens and types, respectively.

The variables in (1b i-ii) count the number of the phonological neighbors (if we apply replacement or replacement, deletion, insertion and transposition, respectively). The variable in (1b iii) is a less strict measure of phonological distance that calculates the mean phonological distance of the N (typically 20) nearest items.

² Detailed information on the nature and the calculation of the variables is available on the ILSP webpage (<http://speech.ilsp.gr/iplr>, see also Protopapas et al. 2010).

The variables chosen allow us to control whether the constructed words are close to but yet not too distant from existing ones. There is, however, a major problem with the Clean Corpus for someone who wants to use it in order to extract results for a morphology-oriented linguistic structure such as Greek stress. Clean does not provide any information on the morphological category (e.g., nouns, verbs, pronouns, etc.) of listed words, which is of absolute relevance for the study at hand, due to the morphology-oriented nature of Greek stress. As argued in Revithiadou (1999), there is a sharp difference between verbs and nouns in their stress behavior; for instance, nouns exhibit more accentual contrasts than verbs. Moreover, in verbs stress is transparently associated to morphological information. For example, past forms are almost exclusively associated with APU stress.³

The solution to this problem was to develop a finer grained, noun-targeted version of Clean, named *NClean Corpus*. *NClean* consists of 13.324 (non-derived/non-compound) nouns, all culled up from the Clean Corpus, version: ignoring stress. We relied on the stressless version of the corpus because, given the aims of our study, we didn't want the variables in (1) to take into consideration in their calculations information on the position of stress. The next step was to extract valuable information contained in *NClean* and effectively exploit it in the construction of pseudowords for our experimental procedures.

³ APU stress is affiliated with the PAST either because past inflections have been (traditionally) argued to require stress to surface on the APU syllable (Warburton 1970, Babiniotis 1972, Ralli 2005) or because a stressed proclitic or prefixal element is present in the past form (see van Oostendorp 2007, 2012 and Spyropoulos & Revithiadou 2009, 2011, respectively).

3.2. Constructing the pseudowords

The 200 pseudowords created for our experimental tasks on noun stress were controlled for morphological classhood, size and syllable structure. More specifically, two- and three-syllable long nouns from the five noun classes: *-os*, *-o*, *-a*, *-as*, and *-i* (fem/neut) were constructed. We opted for simple syllable structures, namely, CV.CV(C), CV.CV.CV(C), CCV.CV(C), CCV.CV.CV(C), in order to avoid the possible interference of phonotactics in our experimental research. Tables 1 and 2 show examples of masculine and feminine nouns, respectively. Since stress is the main quest of our research it not assigned in the pseudowords.

<i>-as</i>	2σ	3σ
CV.CVC	θokas	
CCV.CVC	krefas	
CCV.CV.CVC		trivetas
CV.CV.CVC		lavenas

Table 1: Masculine pseudonouns in *-as*

<i>-a</i>	2σ	3σ
CV.CV	rova	
CCV.CV	spika	
CV.CV.CV		letoma
CCV.CV.CV		krixena

Table 2: Feminine pseudonouns in *-a*

The following procedure was followed for each constructed word:

- STEP 1: All data of the NClean Corpus nouns were categorized according to size and syllable structure. As a result, nouns of the same length and syllable structure were grouped together.

(7) disyllabic CCV.CV nouns

	A	B	C	D	E	F	G
1	spel	phon	BGtokfreqPho	BGtypfreqPho	nNeiPho	nNeiRDITPho	PLD20
396	σκουφι	skufi	0,112	0,351	5	7	1.650
397	σκουφο	skufo	0,116	0,323	5	6	1.600
398	σκουφοι	skufi	0,112	0,351	5	7	1.650
399	σκουφου	skufu	0,068	0,200	4	5	1.700
400	σκυλα	scila	0,703	0,955	10	13	1.350
401	σκυλε	scile	0,711	0,904	7	7	1.600
402	σκυλι	scili	0,902	1,155	15	18	1.100
403	σκυλια	scila	0,147	0,241	8	8	1.650
404	σκυλιου	scilu	0,083	0,133	6	6	1.750
405	σκυλο	scilo	0,745	1,026	13	16	1.150
406	σκυλοι	scili	0,902	1,155	15	18	1.100
407	σκυλου	scilu	0,516	0,688	10	12	1.400

(8) trisyllabic CCV.CV.CVC nouns

	A	B	C	D	E	F	G
1	spel	phon	BGtokfreqPho	BGtypfreqPho	nNeiPho	nNeiRDITPho	PLD20
588	πριγκηπας	prihGipas	0,715	0,833	2	4	2150
589	πριγκηπες	prihGipes	0,703	0,778	2	2	2300
590	πριγκηπων	prihGipon	0,844	0,874	0	0	2600
591	πριγκυπας	prihGipas	0,715	0,833	2	4	2150
592	πριγκυπες	prihGipes	0,703	0,778	2	2	2300
593	πριγκυπων	prihGipon	0,844	0,874	0	0	2600
594	προβατον	provaton	1,511	1,616	0	1	1950
595	προβατων	provaton	1,511	1,616	0	1	1950
596	προβολεις	provolis	1,217	1,487	5	9	1550
597	προβολες	provoles	0,934	1,162	5	7	1650
598	προβολης	provolis	1,217	1,487	5	9	1550
599	προβολος	provolos	1,011	1,282	5	5	1750

- STEP 2: Mean values and SDs for bigram frequencies (phonemes only) and neighborhoods & cohorts were calculated anew for each noun category (e.g, for disyllabic CV.CVC nouns in *-os*, *-as*, disyllabic CCV.CV nouns in *-o*, *-a*, *-i*, etc., trisyllabic CCV.CV.CVC nouns in *-o*, *-a*, *-i*, and so on). We set the acceptable range as strictly as possible, from mean $-1SD$ to mean $+1SD$. For instance, in nouns with syllable structure CV.CV, the mean value of BGtokfreqPho was 1,001 and the SD was 0,866. Thus, the permissible range was set from [mean value - SD = 1,001 - 0,866 =] **0,135** to [mean value + SD

= 1,001 + 0,866 =] **1,867**. Similarly, the mean value of nNeiPho for the same category of nouns was 18 and the SD was 10. Hence the permissible range was set from [18 - 10 =] **8** to [8 + 10 =] **28**.

- STEP 3: Novel words were constructed and tested by the NumTool (<http://speech.ilsp.gr/iplr/NumTool.aspx>, see Protopapas et al. 2010), which provided quantitative measures of the variables in question for each submitted word string.

(9)

NUM Tool		Logmean bigram token frequency (phonemes only)	Logmean bigram type frequency (phonemes only)	N phonological neighbors (standard: replace only)	N phonological neighbors (replace,delete,insert,transpose)	Phonological Levenshtein distance 20	
Enter up to 20 words or nonwords: ζακα κιντα χιπα μπαρος τουζος λαμος		ζακα	0.449	0.926	13	13	1.450
		κιντα	0.408	0.617	13	13	1.450
		χιπα	1.224	1.208	12	13	1.400
		μπαρος	0.943	1.432	16	17	1.300
		τουζος	0.319	0.403	1	1	1.950
		λαμος	1.252	1.621	14	17	1.150

- STEP 4: Words that fell within the defined range of all variables at hand (see Step 2) were selected as suitable items for the experiment. Those that failed to fit to the defined range of at least one variable were discarded as unsuitable.

Some representative examples of pseudowords are provided in Tables 3 and 4. Each table presents the range of each variable within the specific category (feminine *-a* and masculine *-os*, respectively). In order for a constructed pseudoword to be accepted, its values for all variables should be within the appropriate range. Novel words whose values deviated from a given range (e.g., *τουζος* /tuzos/) were excluded as experimental items.

Feminine nouns in -a, 2σ, syllable type: CV.CV						
Pseudowords		BGtok freqPho	BGtyp freqPho	nNei Pho	nNei RDITPho	PLD20
		0,135-1,867	0,305-1,996	8-28	9-33	0,942-1,562
ζακα	zaka	0,449	0,926	13	13	1,450
κιντα	kida	0,408	0,617	13	13	1,450
χιπα	xipa	1,224	1,208	12	13	1,400

Table 3: Feminine pseudonouns in -a

Masculine nouns in -os, 2σ, syllable type: CV.CVC						
Pseudowords		BGtok freqPho	BGtyp freqPho	nNei Pho	nNei RDITPho	PLD20
		0,408-2,213	0,627-2,335	5-20	6-26	1,015-1,757
μπαρος	baros	0,943	1,432	16	17	1,300
τουζος	tu:zɔs	0,319	0,403	1	1	1,950
λαμος	lamos	1,252	1,621	14	17	1,150

Table 4: Masculine pseudonouns in -os

The end result was a pool of words that complied to the defined value range of the respective word categories of the NClean Corpus and their degree of familiarity to the native speakers was well-defined and measurable, as intended.

4. Conclusions

This study demonstrates the usefulness of corpora and the associated quantitative tools in constructing experimental material that complies to the phonotactic restrictions and, in general, to the phonological structure of Greek. More importantly, it establishes a methodology that can weigh rather accurately the familiarity effect a word may have to the native speakers' ears. Finally, it shows that the incorporation of morphological information enhances the applicational power of corpora leading towards more targeted results.

References

- Alderete, John. 1999. *Morphologically governed accent in Optimality Theory*. PhD dissertation, University of Massachusetts, Amherst.
- Alderete, John. 2001a. *Morphologically governed accent in Optimality Theory*. New York: Routledge.
- Alderete, John. 2001b. Dominance effects as transderivational anti-faithfulness. *Phonology* 18: 201-253.
- Apostolouda, Vasso. 2012. *Ο Τόνος των Ουσιαστικών της Ελληνικής: Μια Πειραματική Προσέγγιση*. [Nominal stress in Greek: An experimental approach.] MA thesis, Aristotle University of Thessaloniki.
- Babinotis, Georgios. 1972. *Το Ρήμα της Ελληνικής* [The Verb in Greek]. Athens.
- Burzio, Luigi & Niki Tantalou. 2007. Modern Greek accent and faithfulness constraints in OT. *Lingua* 117: 1080-1124.
- Crosswhite, Katherine, John Alderete, Tim Beasley & Vita Markman. 2003. Morphological effects on default stress placement in Russian novel words: An experimental approach. In Gina Garding & Mimu Tsujimura (eds.) *WCCFL 22 Conference Proceedings*. Somerville, MA: Cascadilla Press, 151-164.
- Drachman, Gaberell and Angeliki Malikouti-Drachman. 1999. Greek word accent. In Harry van der Hulst (ed.) *Word Prosodic Systems in the Languages of Europe*. Berlin & New York: Mouton de Gruyter, 897-945.
- Fainleib, Lena. 2008. *Default stress in unpredictable stress languages: evidence from Russian and Hebrew*. MA dissertation, Tel Aviv University.
- Halle, Morris. 1973. The accentuation of Russian words. *Language* 49: 312-348.
- Halle, Morris. 1997. On stress and accent in Indo-European. *Language* 73: 275-313.
- Kiparsky, Paul & Morris Halle. 1977. Toward a reconstruction of the Indo-European accent. In Larry M. Hyman (ed.) *Studies in Stress and Accent*. Los Angeles: University of Southern California, 209-238.
- Lavitskaya, Yulia and Barış Kabak. 2011a. Russian accentual system revisited: experimental and diachronic evidence. Paper presented at *OCP8* (January 20-23, 2010), University of Hassan II, Ain Chock, Marrakesh, Casablanca.
- Lavitskaya, Yulia and Barış Kabak. 2011b. Default stress in Russian: An experimental study. Paper presented at the *International Workshop on Suprasegmentals in Acquisition and Processing* (May 31-June 01, 2011), University of Konstanz, Konstanz.
- Malikouti-Drachman, Angeliki and Gaberell Drachman. 1989. Stress in Greek. In *Studies in Greek Linguistics 1989*. University of Thessaloniki: 127-143.
- Melvold, Janis Leanne. 1990. *Structure and Stress in the Phonology of Russian*. Doctoral dissertation, MIT, Cambridge, MA.
- Nikolaeva, Tatiana. 1971. Mesto udareniiia i foneticheskii sostav slova (rasstanovka udareniiia v neizvestnykh slovakh inostrannogo proizkhozheniia). In Fedot P. Filin et al. (eds.) *Fonetika. Fonologiya. Grammatika. K semidesiatiletiiu A.A. Reformatskogo* [Phonetics. Phonology. Grammar. For the 70th Birthday of A.A. Reformatskii]. Moscow: Nauka.
- Nikolou, Kalomoira, Anthi Revithiadou & Despina Papadopoulou. 2012. Exceptional stress patterns in the absence of morphological conditioning. Paper presented at the *10th International Conference on Greek Linguistics* (September 1-4, 2011). Komotini: Democritus University of Thrace.
- Oostendorp, Marc van. 2007. Derived Environment Effects and Consistency of Exponence. In: Sylvia Blaho, Patrik Bye & Martin Krämer (eds.), *Freedom of Analysis?*. 123-148, Berlin: Mouton De Gruyter.
- Oostendorp, Marc van. 2012. Stress as a proclitic in Modern Greek. *Lingua* 122: 1165-1181.
- Protopapas, Athanassios, Svetlana Gerakaki & Stella Alexandri. 2006. Lexical and default stress assignment in reading Greek. *Journal of Research in Reading* 29(4): 418-432.
- Protopapas, Athanassios, Tzakosta, Marina, Chalamandaris, Aimilios & Pirros Tsiakoulis. 2010. IPLR: An online resource for Greek word-level and sublexical information. *Language Resources and Evaluation*. Retrieved 27 October 2013 from: <http://link.springer.com/article/10.1007%2Fs10579-010-9130-z>
- Ralli, Angela & Loudovikos Touratzidis. 1992. A computational treatment of stress in Greek inflected forms. *Language and Speech* 35: 435-453.
- Ralli, Angela. 2005. *Μορφολογία* [Morphology]. Athens: Patakis.
- Revithiadou, Anthi. 1999. *Headmost Accent Wins: Head Dominance and Ideal Prosodic Form in Lexical Accent Systems*. Doctoral dissertation, LOT Dissertation Series 15 (HIL/Leiden University), The Hague: Holland Academic Graphics.

- Revithiadou, Anthi. 2007. Colored Turbid accents and Containment: A case study from lexical stress. In Sylvia Blaho, Patrick Bye & Martin Krämer (eds.) *Freedom of Analysis?*. Berlin and New York: Mouton de Gruyter, 149-174.
- Revithiadou, Anthi, Aggelos Lengeris & Dimitra Ioannou. In prep. In search of the default stress in Greek. Ms., AUPh.
- Spyropoulos, Vassilios & Anthi Revithiadou. 2009. The morphology of PAST in Greek. *Studies in Greek Linguistics* 29: 108-122.
- Spyropoulos, Vassilios & Anthi Revithiadou. 2011. PAST in Greek: A case study of the interface between morphological structure and phonological realization. Ms. University of Athens and University of Thessaloniki.
- Topintzi, Nina & Evia Kainada. 2012. Acronyms and the placement of default stress in Greek. Paper presented at the *10th International Conference on Greek Linguistics* (September 1-4, 2011). Komotini: Democritus University of Thrace.
- Warburton, Irene. 1970. *On the Verb of Modern Greek*. The Hague: Mouton and Co.