



## Specific features of the Galician language and implications for speech technology development

Manuel González González, Eduardo Rodríguez Banga, Francisco Campillo Díaz, Francisco Méndez Pazó, Leandro Rodríguez Liñares, Gonzalo Iglesias Iglesias

### ► To cite this version:

Manuel González González, Eduardo Rodríguez Banga, Francisco Campillo Díaz, Francisco Méndez Pazó, Leandro Rodríguez Liñares, et al.. Specific features of the Galician language and implications for speech technology development. *Speech Communication*, 2008, 50 (11-12), pp.874. 10.1016/j.specom.2008.02.004 . hal-00499205

**HAL Id: hal-00499205**

**<https://hal.science/hal-00499205>**

Submitted on 9 Jul 2010

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Accepted Manuscript

Specific features of the Galician language and implications for speech technology development

Manuel González González, Eduardo Rodríguez Banga, Francisco Campillo Díaz, Francisco Méndez Pazó, Leandro Rodríguez Liñares, Gonzalo Iglesias Iglesias

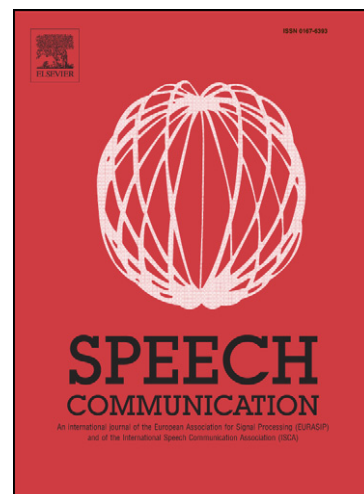
PII: S0167-6393(08)00022-8  
DOI: [10.1016/j.specom.2008.02.004](https://doi.org/10.1016/j.specom.2008.02.004)  
Reference: SPECOM 1690

To appear in: *Speech Communication*

Received Date: 31 May 2007  
Revised Date: 1 November 2007  
Accepted Date: 14 February 2008

Please cite this article as: González, M.G., Banga, E.R., Díaz, F.C., Pazó, F.M., Liñares, L.R., Iglesias, G.I., Specific features of the Galician language and implications for speech technology development, *Speech Communication* (2008), doi: [10.1016/j.specom.2008.02.004](https://doi.org/10.1016/j.specom.2008.02.004)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



# Specific features of the Galician language and implications for speech technology development

Manuel González González<sup>a</sup> Eduardo Rodríguez Banga<sup>b</sup>  
Francisco Campillo Díaz<sup>b</sup> Francisco Méndez Pazó<sup>b</sup>  
Leandro Rodríguez Liñares<sup>c</sup> Gonzalo Iglesias Iglesias<sup>b</sup>

<sup>a</sup>*Dpto. Filoloxía Galega. Universidade de Santiago. Santiago de Compostela.  
SPAIN*

<sup>b</sup>*Dpto. Teoría de la Señal y Comunicaciones. Universidad de Vigo. Vigo. SPAIN*

<sup>c</sup>*Dpto. Informática. Universidad de Vigo. Ourense. SPAIN*

**Corresponding author:** Eduardo Rodríguez Banga.

**Address:**

Dpto. Teoría de la Señal y Comunicaciones. ETSI Telecomunicación.  
Universidad de Vigo. Campus Universitario. 36310. Vigo. SPAIN

**e-mail:** erbanga@gts.tsc.uvigo.es

**Phone:** +34 986 812676 **Fax:** +34 986 812116

---

## Abstract

In this article we present the main linguistic and phonetic features of Galician which need to be considered in the development of speech technology applications for this language. We also describe the solutions adopted in our text-to-speech system, also useful for speech recognition and speech-to-speech translation. On the phonetic plane in particular, the handling of vocal contact and the determination of mid-vowel openness are discussed. On the linguistic plane we place special emphasis on the handling of clitics and verbs. It should be noted that in Galician there is high interrelation between phonetics and grammatical information. Therefore, the task of morphosyntactic disambiguation is also addressed. Moreover, this task is fundamental for a higher level linguistic analysis.

### *Key words:*

Galician, Iberian languages, speech technology, text-to-speech, linguistic analysis

---

## 1 Introduction

Galician is a Romance language — like French, Italian, Romanian, Spanish, Portuguese and Catalan — which resulted from the evolution of Latin in the northwestern part of the Iberian Peninsula.

Galician is mainly spoken in the Autonomous Community of Galicia, but also in the western parts of the provinces of Asturias, León and Zamora, which are located beyond the administrative limits of Galicia. A variety of Galician is also spoken in a small enclave in the northwest of the province of Cáceres, where approximately 5,000 inhabitants preserve the Galician language, brought there by Christian conquerors in the 12th and 13th centuries.

There is no precise information on Galician speaker numbers. According to official statistics, Galicia has a population of about 3 million inhabitants. About 90% of this population can speak fluent Galician, which is the language of daily use for 57% of Galicians, and of occasional use for 30% of Galicians. Nevertheless the total number of speakers is believed to be around 4.5 million, as Galician is also spoken by immigrants to other areas of Spain (especially Barcelona and the Basque Country), to countries in Europe (especially Germany, Switzerland and France), and to countries in South America. For instance, Buenos Aires is said to be the biggest Galician city because of the large number of Galician immigrants living there.

The golden age of the Galician language was during the Middle Ages, when it was the lyric language *par excellence* and used to compose poetry by kings,

troubadours and minstrels in the courts of the Iberian Peninsula and in foreign courts. The Way of St. James undoubtedly played an important role in this period. Use of Galician as a language of culture began to decline by around the end of the 15th century. It eventually became relegated to use as an oral language during the 16th to the 18th centuries. In the 19th century, in a social movement known as the *Rexurdimento* (resurgence), Galician was revived as a language of culture and was claimed by right as the language of Galicia. It was not until 1936, however, that Galician received recognition as an official language in the Statute of Autonomy of Galicia. Unfortunately, the Spanish Civil War prevented the Statute from being implemented. In 1981, with the end of Franco's dictatorship, the Statute of Autonomy of Galicia was finally approved and the Galician language was made an official language of Galicia, along with Spanish.

As a minority language, Galician is, in certain aspects, at a clear disadvantage in relation to other languages such as Spanish or English. This disadvantage extends to the new technologies in general. Since the survival of a language depends on its adaptation to new circumstances and new communication modes, we have focused a great part of our research effort on the development of speech technology for Galician.

Speech recognition, text-to-speech systems and automatic translation systems have evolved enormously in the last decade, with speech-to-speech automatic translation, in particular, gaining increasing importance. In practice a speech-to-speech translation system involves automatic speech recognition (ASR), automatic text-to-text translation and text-to-speech (TTS) synthesis modules. In order to obtain an acceptable translation system, all these modules must be reliable in terms of performance. Performance depends greatly on adequate language modelling, as otherwise errors in the early stages of development will be propagated and magnified in subsequent modules.

Language models used in speech recognition and statistical automatic translation systems are normally based on word n-grams — which represent probabilities for sequences of words. In both cases very little linguistic knowledge is generally considered; rather, text-to-speech systems attempt to make a linguistic analysis of the text to be synthesised, although only superficially in many cases.

In recent years, much of our work has been devoted to improving linguistic analysis in our *Cotovía* text-to-speech system for Galician and Spanish. The automatic determination of the sentence structure and its components is also useful for ASR and statistical translation systems. In these systems, the language models based on word n-grams often generate sentences that are linguistically incorrect (because of lack of gender or number concordance, for example). Taking linguistic information into account would enable incorrect

structures to be detected for further action (penalisation, discarding, etc.)

Obviously, linguistic analysis is highly dependent on individual languages, although there are many similarities among languages with a common origin. All the languages spoken at present in the Iberian Peninsula, with the exception of Basque, come from Latin and so have many common linguistic analysis difficulties.

In this work we focus on the particularities of Galician and related problems of linguistic analysis in some depth. Of the Iberian languages, Spanish is the most widely used and is probably the language in which the greatest research effort has been made. For this reason we will use Spanish as a reference language in order to highlight the distinctive features of Galician. The description of these features is also supported by statistical measures obtained from the automated analysis of a large journalistic corpus – collected from the on-line Spanish and Galician editions of the ‘El Correo Gallego/ O Correo Galego’ between 1997 and 2003 — whose size is about 15 million words for each language.

This article is organised as follows: in Section 2, we describe the main characteristics of Galician phonetics and discuss the impact of contextual factors; in Section 3, we discuss linguistic particularities and certain other phonetic issues related to grammar categories; morphosyntactic analysis is dealt with in Section 4 and in Section 5 we briefly describe intonational aspects. Finally, Section 6 presents our conclusions.

## 2 Galician Phonetics

Although Spanish and Galician depart from a common Latin phonological system, divergent evolution over the centuries has generated significant differences in their vocalic and consonantal subsystems. Several factors have played a central role in this evolutionary process: on the one hand, the natural tendency towards the spatial differentiation of languages over the course of time, and on the other hand, the influence of substratum factors (e.g., the languages spoken in the regions before the arrival of the Romans) and superstratum factors (e.g., the influence of Germanic and Arab invaders).

### 2.1 Vowels

In Galician there are seven vowels, shown in Table 1 and Figure 1 in SAMPA <sup>1</sup> notation, unlike Spanish which has only five (it does not have the oppositions

<sup>1</sup> <http://www.phon.ucl.ac.uk/home/sampa>

[e]/[E] and [o]/[O]). Nevertheless, the two languages use only five graphemes to represent their vowels. In Spanish the automatic phonetic transcription of vowels is very simple, since each grapheme corresponds to a unique phoneme. In Galician, however, it is very difficult to determine mid-vowel openness (graphemes *e*, *o*) because this essentially depends on etymology.

Despite the complexity of determining mid-vowel openness, several studies (Castro (1998), Veiga (1976)) agree on certain points that simplify this task:

- There is at most one open mid-vowel per word.
- Open mid-vowels can only appear in stressed syllables, or in unstressed syllables in the word-initial position. As an exception to this norm, there are some derived forms such as diminutives where the openness degree of the vowels is maintained even though the stress changes to the vowel *i* in the appended suffix (*-iño(s)* or *-iña(s)*).
- Open middle vowels do not appear after stressed syllables, because the oppositions [e]/[E] and [o]/[O] are neutralised in this context.

From a practical point of view, the problem of determining mid-vowel openness mainly affects nouns, verbs and adjectives. With verbs, the problem is even more complex because mid-vowel openness in a verbal root will change depending on the tense, mode, number or person. This phenomenon led us to make an exhaustive analysis of Galician verbal forms, described in Section 3.5 of this paper. The analysis for adjectives and nouns is described in Section 3.1.

In Galician, vocalic duration plays no distinctive role, as long and short vowels in classical Latin were transformed into different grades of aperture in vulgar Latin, the antecedent of Galician.

SAMPA	Grapheme	Description	Example	Transcription
a	a	open central	casa	''ka-sa
E	e	open-mid front unrounded	medo	''mE-Do
e	e	close-mid front unrounded	pena	''pe-na
i	i	close front unrounded	illa	''i-Za
O	o	open-mid back rounded	home	''O-me
o	o	close-mid back rounded	poso	''po-so
u	u	close back rounded	curto	''kur-to

Table 1

Galician vowels (the transcriptions include the syllable boundary and the stress markers).

Another characteristic of the Galician vowel system is that, like Spanish, the grouping of two or three vowels in a single syllable is very frequent (about

8% of the syllables in both languages). Diphthongs can be formed either by a strong vowel (a, e, o) with a weak vowel (i, u) in any order, or by two weak vowels. Triphthongs are necessarily formed by a strong vowel between two weak vowels. In diphthongs and triphthongs the weak vowels are pronounced as semivowels or glides, whose phonetic symbols are represented in Table 2.

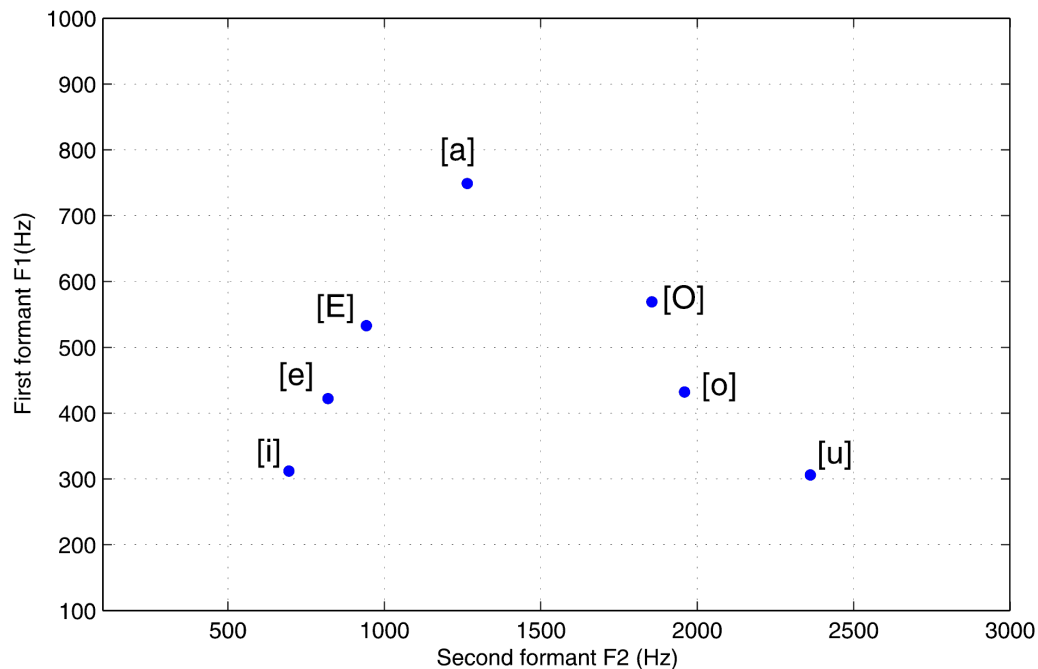


Fig. 1. Formant space of Galician vowels (from (González and Regueira (1994))).

SAMPA	Grapheme	Description	Example	Transcription
j	i	palatal approximant	loito	”loj-to
w	u	labial-velar approximant	fraude	”fraw-De

Table 2

Galician semivowels.

As for the diphthongs, there are certain quantitative significant differences between Spanish and Galician, as follows:

- The frequency of increasing diphthongs in Spanish (about 83% of the diphthongs) is higher than in Galician (about 52%). This is a consequence of the diphthongization process of the vulgar Latin open [E], as in *ie* (PETRA >> *piedra*), and the vulgar Latin open [O], as in *ue* (NOVA >> *nueva*). This process did not take place in Galician (PETRA >> *pedra*, NOVA >> *nova*).
- The frequency of decreasing diphthongs in Galician (about 48%) is higher than in Spanish (about 17%). This is due to the monophthongization process of these diphthongs in Spanish (Gal. *maneira* / Sp. *manera*, Gal. *ouro* / Sp. *oro*).



- c) Spanish has homorganic diphthongs (iu [ju], ui [wi]) that are practically unknown in Galician, since Galician tends to pronounce these sequences as hiatuses (Gal. *mi-ú-do*, *cons-tru-ír*).

Phoneme	Galician	Spanish
a	12.87	12.46
e/E	12.79	13.64
i/j	7.38	7.77
o/O	11.13	9.53
u/w	3.15	3.15

Table 3

Relative frequencies (%) of Galician and Spanish vowels in the journalistic corpus.

Table 3 shows the relative frequencies of the Galician and Spanish vowels, estimated from the phonetic transcription given by our TTS. For the sake of simplicity, weak vowel estimates include the corresponding semivowels. Since open mid-vowels are only considered in stressed syllables (as will be discussed below), their frequency estimates are not accurate. We therefore decided to group the results of open and close mid-vowels in order to obtain a reliable estimate. The estimated frequencies for the Spanish vowels are comparable to those described in (Moreno et al. (2006)) for a written corpus and comparable results were also obtained for Spanish consonants (see next subsection).

## 2.2 Consonants

Table 4 is a classification of Galician consonantal phonemes, and Table 5 describes a set of rules for phonetic transcription that covers the vast majority of cases (excluding some exceptions and foreign words). The set of rules are represented using a notation that is similar to that used for regular expressions, which means that many particular cases can be described in a compact and efficient way.

In Table 5 “#” represents a pause, “\_” is a blank in the text (meaning the beginning or the end of a word) and “\_\*” means zero or more blanks (this allows rules to be applied to intra-word and inter-word grapheme sequences). As an example, let us consider the rule:

$$[\#mn]_{-}^{*}\mathbf{g} \rightarrow \mathbf{g}$$

The string to the left of the arrow, which represents a sequence of graphemes, is the condition of the transcription rule for the grapheme(s) in bold. The right side of the rule denotes the phonetic transcription of the corresponding

grapheme(s) once the left part is verified. More concretely, this rule means that the grapheme **g** is transcribed as [g] when preceded by a pause or by the graphemes *m* or *n*, no matter whether they belong to the same word as **g** or to the previous word. Similarly, the rule

$$\mathbf{n}_{-}[\text{aeiouh}] \rightarrow \text{N}$$

means that every word-final grapheme **n** followed by a vowel or the grapheme *h* is transcribed as the phoneme [N].

The case of the two graphemes **gu** is considered in the following pair of rules:

$$\begin{aligned} \mathbf{gu}[\text{ei}] &\rightarrow \text{g|G} \\ \mathbf{gü}[\text{ei}] &\rightarrow (\text{g|G}) \text{ w} \end{aligned}$$

In these rules “|” means choosing one of two possible phonemes depending on the transcription rules for the grapheme *g*.

It should be emphasised that the application order of the transcription rules is relevant, since rules that affect a greater number of graphemes must be applied first.

With the help of the phonetic transcription rules listed in Table 5 and of the relative frequencies of Galician and Spanish consonants, given in Table 6, we will now discuss certain important features of Galician phonemes.

### 2.2.1 Plosives

In Galician, as in many other languages, there are three unvoiced occlusive phonemes ([p], [t], [k]) and three voiced occlusive phonemes ([b], [d], [g]). While the former may appear in any phonetic context, the latter are often transformed to their approximant variant [B], [D], [G]. For example, transcription rules for the graphemes *b* and *v* (in Galician both have the same sound) state that the phoneme [b] appears after a pause or a nasal (graphemes *m* or *n*). Any other case gives the more frequent approximant variant (as shown in Table 6).

### 2.2.2 Fricatives and Affricates

In Galician, there are five fricatives ([f], [s], [S], [T], [Z]) and one affricate ([tS]). The fact that Spanish does not have [S] means that this phoneme is one of the most characteristic sounds of Galician.

The phonemes [s] and [S] differ in relation to a slight variation at the point

SAMPA	Description	Example	Transcription
p	voiceless bilabial plosive	pato	''pa-to
t	voiceless alveolar plosive	tomo	''to-mo
k	voiceless velar plosive	canto	''kan-to
b	voiced bilabial plosive	bico	''bi-ko
d	voiced dental plosive	dous	''dows
g	voiced velar plosive	gato	''ga-to
B	voiced bilabial approximant	sobre	''so-Bre
D	voiced dental approximant	dedo	''de-Do
G	voiced velar approximant	amigo	a-''mi-Go
f	voiceless labiodental fricative	feira	''fej-ra
s	voiceless alveolar fricative	saco	''sa-ko
S	voiceless postalveolar fricative	xunta	''Sun-ta
T	voiceless dental fricative	berce	''ber-Te
tS	voiceless postalveolar affricate	cheo	''tSe-o
m	voiced bilabial nasal	máis	''maj-s
n	voiced alveolar nasal	nome	''no-me
J	voiced palatal nasal	viño	''bi-Jo
N	voiced velar nasal	unha	''uN-a
l	voiced alveolar lateral	alto	''al-to
Z	voiced postalveolar fricative	fillo	''fi-Zo
r	alveolar flap	paro	''pa-ro
rr	alveolar trill	ría	''rri-a

Table 4  
Galician consonants.

of articulation, which is alveolar in the former and postalveolar in the latter. This variation means that most of the energy for the phoneme [S] occurs in middle frequencies, whereas [s] usually concentrates most of its energy in higher frequencies. Moreover, the second formant frequency is normally higher and the second and third formants are closer in [S] than in [s] (as described in (Olive et al. (1993)) for English [S]). In Figure 2 we represent a segment of the waveform and spectrogram of the word *masaxista*. We can clearly observe the differences between the [S], in the middle, and the two alveolars [s].

The affricate [tS] is very similar to the Spanish but could intuitively seem more

Grapheme(s)	Transcription rules	Examples
p, t, k	<b>p</b> → p, <b>t</b> → t, <b>k</b> → k	pato → "pa-to
b, v	[#mn]_* <b>[bv]</b> → b otherwise → B	ambos → "am-bos ovella → o-"Be-Za
d	[#lmn]_* <b>d</b> → d otherwise → D	onde → "on-de orde → "or-De
g	[#mn]_* <b>g</b> → g otherwise → G	algo → "al-go avogado → a-Bo-"Ga-Do
f, s, z, ch	<b>f</b> → f, <b>s</b> → s, <b>z</b> → T, <b>ch</b> → tS	chove → "tSO-Be
x	<b>x</b> → S (general rule) <b>x</b> → ks (exceptions)	xeito → "Sej-to exame → ek-"sa-me
c	<b>c</b> [ei] → T otherwise → k	cento → "Ten-to caso → "ka-so
m, ñ, ll, l	<b>m</b> → m, <b>ñ</b> → J, <b>ll</b> → Z, <b>l</b> → l	maña → "ma-Ja
n	<b>n</b> _* <b>[bpv]</b> → m <b>n</b> _* <b>[gmnlrscfkqzxd]</b> → N <b>n</b> _[aeiouh] → N <b>n</b> # → N otherwise → n	un bico → "um-"bi-ko lingua → "liN-gwa un home → "uN-"O-me fun # → "fuN non # → "noN
h	<b>nh</b> → N	unha → "uN-a
rr	<b>rr</b> → rr	carro → "ka-rro
r	[#ln_] <b>r</b> → rr otherwise → r	enredo → eN-"rre-Do cara → "ka-ra
gu	<b>gu</b> [ei] → G  g <b>gü</b> [ei] → (G  g)w	anguía → aN-"gi-a alguén → al-"GeN
-gn, -mn, -ps	<b>-gn</b> → n, <b>-mn</b> → n, <b>-ps</b> → s	gnomo → "no-mo
h	<b>h</b> → (no phoneme)	home → "O-me

Table 5

Basic transcription rules for Galician language (see text for details).

frequent in Galician due to a different evolution of several Latin consonantal groups. More concretely, the Spanish [tS] comes fundamentally from the Latin groups CT (OCTO >> Sp. *ocho*, Gal. *oito*), ULT (MULTU >> Sp. *mucho*, Gal. *moito*), and from the groups PL and CL preceded by consonant (AMPLU >>

Phoneme	Galician	Spanish
p	2.77	2.80
t	5.12	4.86
k	4.11	3.96
b	0.27	0.29
d	1.04	1.05
g	0.13	0.11
B	2.02	1.94
D	4.23	3.95
G	0.95	0.94
f	0.90	0.76
s	7.76	7.23
S	0.64	—

Phoneme	Galician	Spanish
T	2.20	2.25
tS	0.12	0.21
m	3.06	2.82
n	4.82	6.54
J	0.26	0.23
N	2.77	0.48
l	2.80	5.23
Z	0.32	0.52
r	5.83	5.89
rr	0.71	0.69
x	—	0.61

Table 6

Relative frequencies (%) of Galician and Spanish consonants estimated from the journalistic corpus, and including the Spanish voiceless velar fricative phoneme [x].

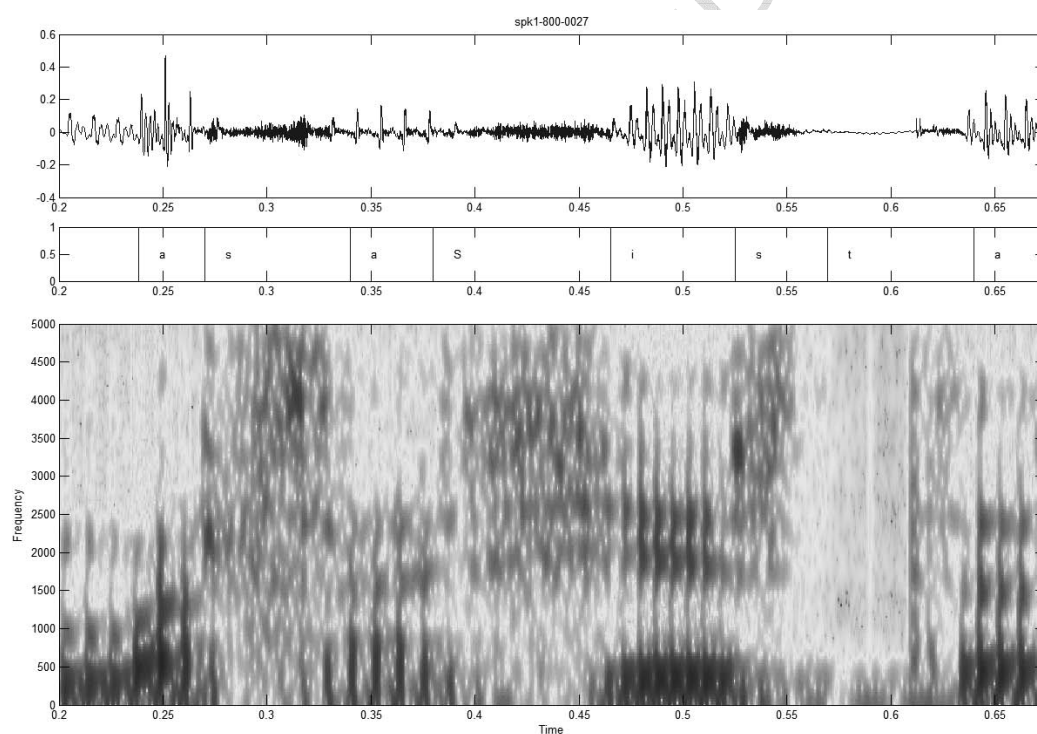


Fig. 2. Waveform and wideband spectrogram of a segment of the word *masaxista*, [ma-sa-<sup>h</sup>Sis-ta], illustrating the differences between the postalveolar [S] and the alveolar [s]

Sp. *ancho*, Gal. *ancho*); while in Galician it comes from the groups PL, CL and FL preceded by a consonant and also in word initial position (PLANU

» Gal. *chan*, Sp. *llano*; CLAMARE » Gal. *chamar*, Sp. *llamar*; FLAMMA » Gal. *chama*, Sp. *llama*). Moreover, in Galician the phoneme [tS] also comes from the group -STI and, therefore, appears in the ending of the second-person singular form of the past simple of verbs (*cantaches*, *comiches*, *fixeches*...). It is also present in the unstressed second-person singular form of the personal pronoun *che*, which is very frequent in spoken Galician.

However, an analysis of the journalistic corpus revealed that, in this case, the frequency of the phoneme [tS] in Spanish is almost twice that for Galician. Nevertheless, we believe that these results are heavily dependent on the style of journalistic texts, generally written in the third person and using a large quantity of numerical data (in Spanish the phonetic transcription of cardinal numerals with the digit 8 produces the phoneme [tS]).

### 2.2.3 Nasals

There are four nasal phonemes: [m], [n], [J], [N]. The phoneme [J], which is shared by most of the Romance languages, is very frequent in Galician not only in the large set of words containing this phoneme, but especially because of its presence in diminutive morphemes, namely, the suffixes *-iño(s)*, *-iña(s)*. Galician makes intensive use of these diminutive forms, especially when spoken.

Undoubtedly, however, one of the most typical sounds of Galician — as much for frequency of use as for use in contexts that are different from other languages — is the velar n, [N]. This sound also appears in other languages (e.g., Spanish, Catalan, English) but always as an allophonic variation of the phoneme [n] before a velar sound ([k], [g]). In Galician it also appears in other consonantal contexts, as reflected in the transcription rules listed in Table 5. Nevertheless, the most striking context of [N] is the intervocalic position. This intervocalic [N] appears as direct transcription of the grapheme group *nh* and whenever a word ends in *n* and the following word begins with a vowel (with just a few exceptions). This latter case is represented by the rule:

$$\mathbf{n}_{-}[\text{aeiouh}] \rightarrow \mathbf{N}$$

The grapheme *n* is also transcribed as [N] in the word final position before a pause. As mentioned, there are a few exceptions to this general rule. More concretely, the final *n* of some pronouns, prepositions and adverbs is transcribed as an alveolar [n] (and not [N]) when followed by an unstressed form of the personal pronoun *o(s)*, *a(s)*. Examples include *quen o fixo*, *alguén a viu*, *ninguén o sabe*, *nín a quixo escoitar*, *sen o mirar*, *ben o sabe el*. This illustrates the importance of morphosyntactic analysis for correct phonetic transcription.

In our Galician TTS about 38% of the graphemes *n* are transcribed as [N]. This proportion is also reflected in Table 6, where the relative frequencies of [N] in Galician and Spanish can also be compared.

#### 2.2.4 *Laterals and Vibrants*

The lateral phoneme [l] and the vibrants [r] and [rr] are the same as in Spanish. The phoneme [l] is much less frequent in Galician than in Spanish, as indicated in Table 6. This is mainly explained by the loss of the intervocalic *l* in Galician patrimonial words (Lat. PALUMBA, Gal. *pomba*, Sp. *paloma*; Lat. CAELU, Gal. *ceo*, Sp. *cielo*; Lat. ANIMALES, Gal. *animais*, Sp. *animales*). With reference to the vibrants, in both Spanish and Galician, [rr] is only possible in the word-initial position or after *l* or *n*, while both [r] and [rr] can appear in the intervocalic position. In every other context the opposition of these phonemes is neutralised.

### 2.3 *The influence of contextual factors*

Contextual factors have a great influence in the speech. Therefore, they should be considered in TTS and ASR systems. This section describes some important and predictable contextual effects in Galician.

#### 2.3.1 *Sound relaxation at the end of elocution*

In (Castro (1998)) the author emphasises the high grade of muscular tension required to pronounce the five vowels in Spanish, even in the unstressed position. She contrasts this with Galician, where only the open mid-vowels are considered to be tense. Our experience shows that these differences between Galician and Spanish are especially noticeable in the final portion of an elocution. Moreover, this relaxation mainly takes place after the last stressed syllable and involves all the phonemes and not just vowels. The same relaxation also occurs in Spanish, but is much less marked than in Galician. As an example, Figure 3 show the waveform and the spectrogram for the final segment of a Spanish sentence pronounced by a Castilian speaker, whereas Figure 4 represents the same sentence in Galician pronounced by a Galician speaker. The phonetic content is practically identical in both cases, which facilitates the comparison. Clearly observable in both cases is the different relaxation that takes place after the last stressed syllable.

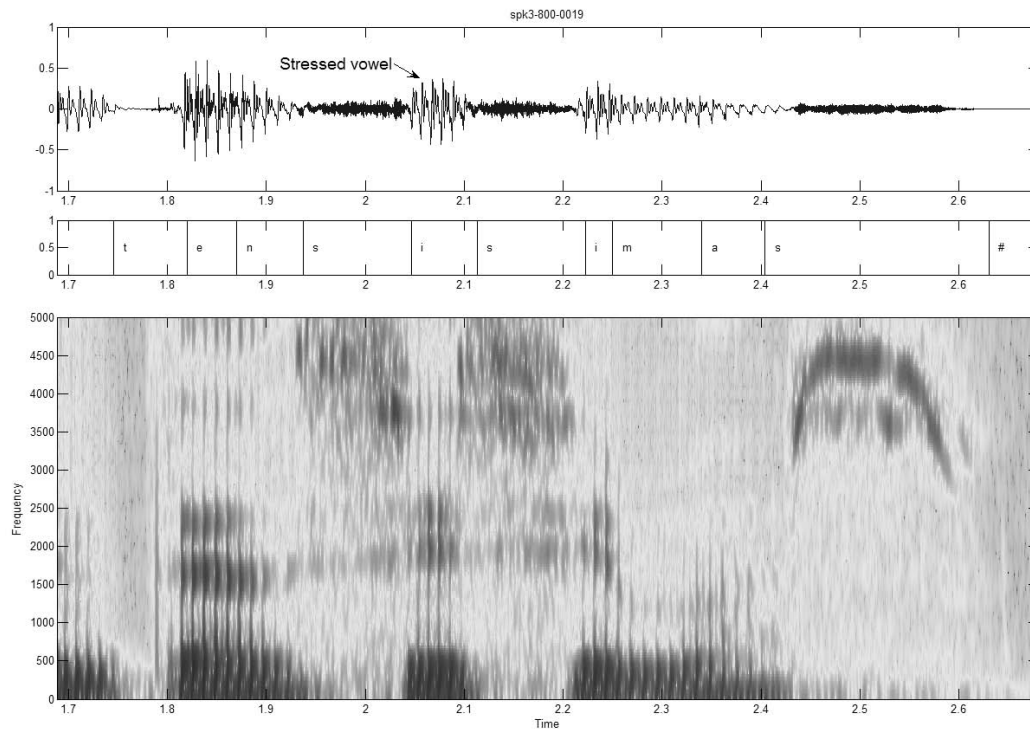


Fig. 3. Waveform and wideband spectrogram of the final segment of the Spanish sentence “Las lluvias eran intensísimas” (see text for details).

### 2.3.2 Vocalic elisions and assimilations

In spoken Galician, vocalic elisions and assimilations are very frequent — much more so than in Spanish. Although present in the written form, many word-final vowels are omitted or transformed when followed by a word beginning with a vowel. The following cases of elision of the word-final vowel are especially common:

- a) Prepositions *ante*, *bardante*, *conforme*, *consonte*, *de*, *dende*, *desde*, *durante*, *entre*, *perante*. Examples: *de agora* [da-’Go-ra] instead of [de-a-’Go-ra], *entre amigos* [en-tra-’mi-Gos] in place of [en-tre-a-’mi-Gos], *desde América* [’des-Da-’mE-ri-ka] instead of [’des-De-a-’mE-ri-ka].
- b) The pronoun or conjunction *que* and the conjunction *porque*. Examples: *o lobo que ouvea* [o-’lo-Bo-kow-’Be-a] instead of [o-’lo-Bo-ke-ow-’Be-a], *dixo que había fame* [’di-So-ka-’Bi-a-’fa-me] in place of [’di-So-ke-a-’Bi-a-’fa-me].
- c) The demonstratives *este*, *ese* : *este amigo* [es-ta-’mi-Go].
- d) The numerals *vinte*, *trinta*, *corenta...*, *noventa*, followed by the conjunction *e*. Examples: *vinte e sete* [’Bin-tE-’sE-te], *trinta e dúas* [’trin-tE-’Du-as], *oitenta e catro* [oj-’tEn-tE-’ka-tro].
- e) The conjunction *se*. Examples: *preguntoulle se había vir axiña* [pre-Gun-’tow-Ze-sa-’Bi-a-’Bi-ra-’Si-Ja].
- f) The adverbs *case*, *onde*, *lonxe*, *sempre*, *tarde*, *hoxe*, *onte*, *antonte*, *de-*



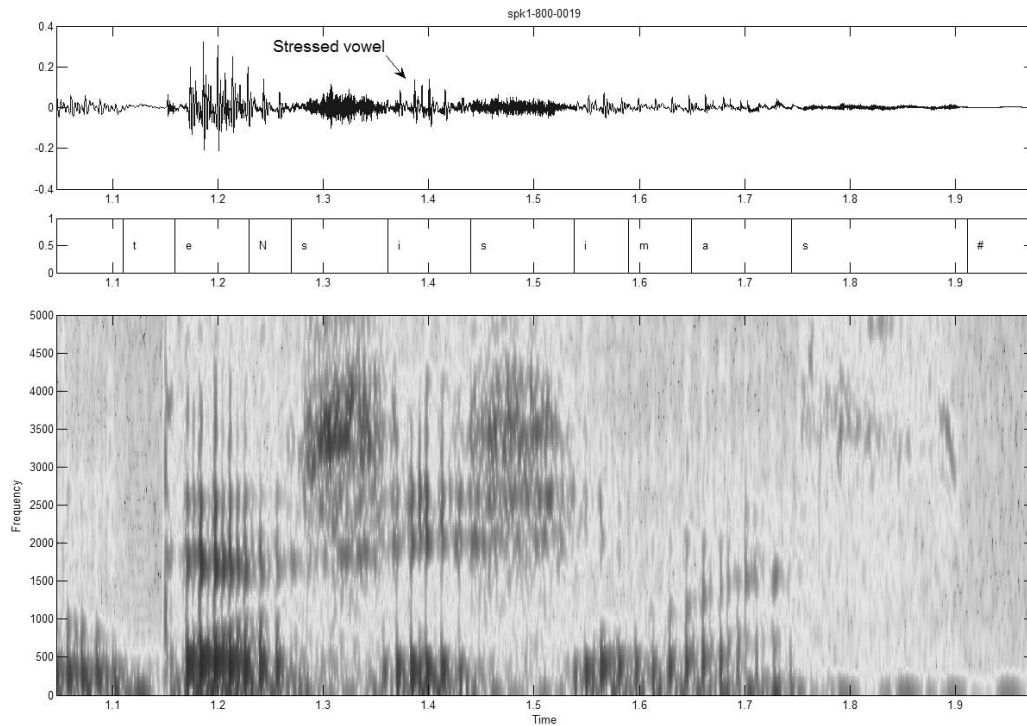


Fig. 4. Waveform and wideband spectrogram of the final segment of the Galician sentence “As choivas eran intensísimas” (see text for details).

*cote*, and adverbs ending in *-mente*. Examples: *sempre andou con contos* [’sem-pran-’dow-koN-’kon-tos], *unanimemente aceptado* [u-na-ni-me-’men-ta-TEp-’ta-Do].

g) The pronouns *me*, *te*, *che*, *lle*, *se*. Example: *deixoume alí* [dej-’Sow-ma-’li].

Some encounters between vowels *a* and *o* are special cases, as follows:

- a) An encounter between vowels *a* and *o* in the post-tonic position is generally resolved as a reciprocal assimilation, whose result is [O]. Example: *fixérao* [fi-’Se-rO].
- b) The same tendency is observed between the unstressed word-final vowel *a* and the unstressed word-initial *o*. Example: *tiña olleiras* [’ti-JO-’Ze-jras]. Nevertheless the assimilation is systematic when a preposition (*ata*, *contra*, *canda*, *deica*, *fóra*, *malia*, *onda*, *para*, *xunta*) or conjunction (*ca*, *coma*) ending in *a* is followed by the article *o(s)*, and when an adverb ending in *a* (*agora*, *axiña*, *inda*, *nunca*, *outrora*, *noutroa*, *de sobra*, *nada*, *seica*) is followed by the article *o(s)* or the unstressed personal pronoun *o(s)*. Examples: *mellor ca o teu* [me-’Zor-kO-’tew], *onda o neno* [’on-dO-’ne-no], *nunca o vin* [’nuN-kO-’BiN].

Finally, we should mention elisions and assimilations that occur between word-final *a* and word-initial *e* when both are unstressed:

- a) An unstressed word-final *a* followed by unstressed word-initial *e* which forms a syllable with one of the consonants *n*, *m*, *l*, *s*: the initial *e* usually disappears and the final *a* forms a syllable with the following consonant. Example: *a miña enfermidade* [a-''mi-Jan-fer-mi-''Da-De], *nesa esquina* [ne-sas-''ki-na].
- b) An unstressed word-final *a* followed by an unstressed word-initial *e* which itself forms a syllable: a reciprocal assimilation occurs whose phonetic result is generally [E]. Example: *mala educación* [''ma-lE-Du-ka-''TjoN].

### 3 Category-Specific Issues

This section discusses certain phonetic and morphological issues that depend on the grammatical category of the word. It is thus organised in terms of several subsections that correspond to the different categories under discussion.

#### 3.1 Nouns and Adjectives

We mentioned above that one of the most difficult problems with nouns and adjectives is to determine mid-vowel openness. This task is especially important for vowels in stressed syllables (as discussed in Section 2.1). Although mid-vowel openness fundamentally depends on etymology, the study of a large number of cases makes it possible to define a set of rules based on sets of terminations and exceptions. In *Cotovía* we use rules to determine stressed mid-vowel openness. The algorithm is mainly based on the position of the stressed syllable within the word (noun or adjective), as follows:

- (1) **Proparoxytone word.** Open by default, with the exception of a termination list (-*ébeda*, -*céntrico*, -*espera*, ...) and a set of words (*alvéolo*, *bébedo*,...)
- (2) **Paroxytone word.** Open by default, except for vowels in contact with a nasal phoneme (with some exceptions), a long list of terminations (-*boca*, -*cerca*, -*eba*, ...) and a set of exceptions.
- (3) **Oxytone word.** Close by default, except for terminations -*é*, -*el*, -*én*, -*en*, -*ol*, -*oz*.

In addition to applying the previous rules, diminutives, which are very frequent in Galician, need to be detected. In general, diminutives are formed by appending the suffixes -*inho(s)*, -*inha(s)*, although often a slight modification of the noun termination is necessary (examples include *can* → *cancinho*, *nenó* → *neniño*). The problem with the diminutives is that the stressed vowel is now the *i* of the suffix, although the mid-vowel openness is the same as in

the original noun. Therefore, in terms of text analysis, we need to be aware that any word ending in *-iño(s)* or *-iña(s)* may be a diminutive, making the determination of the mid-vowel openness a little more complex.

Obviously the accuracy of these rules depends on the exception lists. Nevertheless, as a guideline, in the current version of our TTS accuracy was 92.82% evaluated over a text of about 8,000 words revised by two linguists.

### 3.2 *The second form of the article*

The forms of the definite article (*o, a, os, as*) can be transformed into (*-lo, -la, -los, -las*) under certain circumstances. This transformation is already implicitly considered in many Galician contractions (examples: *por + o* → *polo*, *ambas + as* → *ámbalas*), but in other cases the use of the second form of the article is optional in the written language. Thus, it can be employed when the definite article follows a verbal form that ends in *-r* or *-s* (possible clitics included), with the verb losing the final *-r* or *-s*. As an example, the sentences “*Comer o caldo*” and “*Come-lo caldo*” are both correct and equivalent in written Galician. Nevertheless the recommended pronunciation is the same in both cases and corresponds to the phonetic transcription of the second option.

### 3.3 *Contractions*

We have already referred to the frequency and importance of elisions in Galician. As the language evolved, elisions gave rise to a large number of contractions in written form. This phenomenon contrasts enormously with Spanish, in which only the contractions *al* (*a+el*) and *del* (*de+el*) exist, both of the type preposition+article. Table 7 shows the most common types of Galician contractions and the number of contractions of each type. As can be observed, the number of possible contractions is very large. In addition, the prepositions that form contractions are those most frequently used in Galician. Indicating the importance of this reduced group of prepositions is Figure 5; the histogram on the left shows the most frequent prepositions in the journalistic corpus, irrespective of whether they were found as a part of a contraction or in isolated form, whereas the histogram on the right shows that around 50% of prepositions appeared in contracted form.

Type of contraction	Cases	Example
Preposition (a, en, de, con, por, tras) + Definite Article	24	do (de+o)
Preposition (en, de, con) + Indefinite Article	12	cunha (con+unha)
Preposition (en, de) + 3 <sup>rd</sup> person Personal Pronoun	8	nel (en+el)
Preposition (en, de) + Demonstrative	30	naquel (en+aquel)
Preposition (en,de)+ Indefinite (algún, outro)	16	nalgún (en+algún)
Demonstrative + Indefinite (outro)	12	estoutro (este+outro)
Preposition (en, de) + Demonstrative + Indefinite (outro)	24	destoutro (de+este+outro)
Conjunction (ca, mais) + Definite Article	8	có (ca+o)

Table 7

Main Galician contractions and number of contractions of each type.

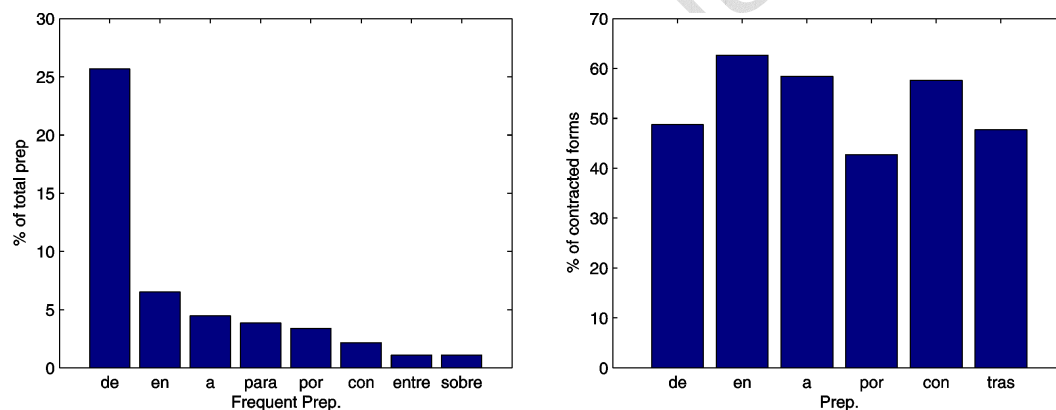


Fig. 5. Most frequent prepositions (left) and percentage of contracted forms (right).

### 3.4 Clitics

In Galician, as in Spanish, the unstressed forms of personal pronouns can go in the proclitic or enclitic position in relation to the verb (before or after the verb, respectively). They are almost equally frequent in the two languages: about 20-25% of the verbal forms of the analysed corpora had clitics. Nevertheless, there are some very noticeable differences, as follows:

- Enclitic forms are much more frequent in Galician than in Spanish, and the enclitic position is more common in Galician than the proclitic position. As an example, in the journalistic Galician corpus about 61% of the unstressed

forms of personal pronouns were in the enclitic position, while the remaining 39% were in the proclitic position (17% and 83% in Spanish, respectively). In Galician, in contrast to Spanish, the third person pronouns in the accusative case *o(s)/lo(s)* and *a(s)/la(s)* in the enclitic position lead to verbal forms ending in *-r* or *-s* losing this last character.

- The sequence of enclitics attached to the verb is usually longer in Galician than in Spanish. In Galician sequences of two and even three pronouns are very usual. In Spanish the use of enclitics is less frequent, and when enclitics are used, only one clitic is appended to the verb, and more rarely, two.
- In Galician clitics (proclitics or enclitics) are often combined to form conglomerates of unstressed pronouns. The unstressed personal pronouns and their possible conglomerates make a total of about 70 clitic sequences. Examples: *non **cho** dixeron* (*che+o* → *cho*), *dixéron**nolo*** (*nos+o* → *nolo*).
- Galician has a clitic that is unknown in Spanish: the solidarity pronoun. Its aim is to make the listener partake in the speaker's attitude. This pronoun formally coincides with the dative case of the second and third person pronouns. Examples: *é**ch**eme moi difícil falar dos meus pais*; *señor doutor, dó**lle**me moito esta perna*.

### 3.5 Verbs

We previously discussed particular features of Galician verbal forms. On the phonetic plane we mentioned the difficulty with determining the openness of mid-vowels in the verbal root. On the morphological plane we referred to clitics and the second form of the definite article. These factors make the analysis or inflection of Galician verbs much more complex than in Spanish.

From a morphological point of view, Galician verbs are usually classified into three classes according to the ending of their infinitive: verbs ending in *-ar*, *-er* and *-ir*. However, there are some verbs whose infinitive does not fall into these classes (*pór*, *compór*, *dispór*, for instance). Without pretending to make an in-depth morphological description of Galician verbs, certain characteristics must be emphasised:

- a) In Galician, there are no compound tenses with auxiliary verbs (*haber*) in contrast to Spanish (Gal. *fixera*, Sp. *había hecho*).
- b) Galician infinitives have personal and impersonal forms while Spanish infinitives only have an impersonal form.
- c) An alternating phenomenon in the openness of the mid-vowels in the stem of the verbs ending in *-er* (*b[e]bo/ b[E]bes... c[o]mo/c[O]mes...*).

There are also other morphological differences with Spanish verbs (different stress positions in some verbal forms, differences in endings, etc.) but they are

not immediately relevant to the purpose of this article.

In Laverca<sup>2</sup> (González González et al. (2002)) we proposed a set of models or paradigms for Galician verb inflection. These models were developed taking into account both written-form changes and phonetic variations. Thus, some verbs traditionally considered to belong to the same inflection model were established as different models for phonetic reasons. The outcome was a set of 116 paradigms that enable any Galician verb to be analysed or inflected.

Since the resulting paradigms take into account the openness of stressed mid-vowels, in our TTS the phonetic transcription of verbal forms is generally free of error. Nevertheless, errors in part-of-speech tagging may induce errors in the phonetic transcription of verbs that were not recognised as such by the morphosyntactic tagger.

We will now describe the main characteristics and implementation of the verb analyser that we are currently using in our text-to-speech system and for automatic translation research (in this last case together with an inflection module).

The verbal analysis and inflection modules are based on two information tables: one consisting of verbal stems and the other consisting of inflectional endings. The structure of both tables is shown in Figure 6.

Stems	Endings
...	...
0,vir,58	fora,(67,48),(69,48),(61,54),(63,54)
abacel,abacelar,1	íra,(67,31),(69,31)
abaf,abafar,1	sa,(184,89),(188,89),(192,89)
abafall,abafallar,1	ta,(184,88),(188,88),(192,88)
...	ita,(192,38)
abastec,abastecer,9	eita,(188,19)
abastez,abastecer,10	ida,(188,9),(188,17),...,(188,87)
abat,abater,8,87	era,(61,8),(63,8),(61,9),(63,9),...
....	...

Fig. 6. Tables of stems and endings (see text for details).

Each row of the stem table contains the verbal stem, its corresponding infinitive and a series of numeric identifiers that represent the sub-models assigned to each stem-infinitive pair. Each verbal model is composed of one or more sub-models, so that a sub-model often covers only a part of a verbal paradigm. For instance, if we focus on the infinitive *abastecer*, which appears

<sup>2</sup> Laverca is a dictionary of Galician verbs with a detailed description of the different inflection models and a software application that allow analysing, inflecting and reading with synthetic speech any Galician verb.

twice in the stem table in Figure 6, we can see that it is conjugated with the stem *abastec-* according to sub-model 9 and with the stem *abastez-* according to sub-model 10.

Each row of the ending table specifies an ending and a sequence of numeric pairs. The second number in each pair denotes the sub-model, whereas the first number in the pair is a unique identifier of the mood, tense, number and person. If we consider the ending table in Figure 6 and look for pairs containing sub-model 9, we find:

ida,(188,9),(188,17),...,(188,87)  
era,(61,8),(63,8),(61,9),(63,9),...

Combining this information with the stem *abastec-* we obtain the verbal forms *abastecida* for code 188 (feminine singular participle) and *abastecera* for codes 61 and 63 (first and third persons of singular of the past perfect of indicative).

The analysis process for a verbal form is more complex: first we look for all the possible endings of a verbal form, and then we search for stems with an inflection sub-model in common with the previously obtained endings. Obviously, this process can throw up several stems, which indicates that the analysed form belongs to several infinitives. This analysis process is further complicated by the frequent presence of enclitics and the second form of the article attached to the verbal form. Three examples will clarify the procedure.

Figure 7 depicts an analysis of the verbal form *Fose*. First we check whether the verbal form contains the second form of the article. As this is not the case, we extract all the possible enclitics. This stage results in two possibilities: with the enclitic *-se* and with no enclitic. Once we have discarded the enclitics, we compare the terminations of the two alternative forms with those in the ending table. We find three potential endings, for which we obtain their corresponding stems. As can be observed in Figure 7 the stem can be empty ( $\emptyset$  in the stem table). This figure also shows the relevant parts of the stem and ending tables for this example. One of the possibilities found for *Fose* is the empty stem followed by the ending *-fose*. Looking at the stem table, the empty stem may correspond to several infinitives (far more than those shown in Figure 7). Comparing the sub-model codes associated with each stem-infinitive pair with those for the ending *-fose* we find that codes 48 and 54 are common; consequently, the verbal form *Fose* comes from the infinitives *ir* or *ser*. The information about tense, mood, number and person is obtained from the first number in each number pair, with sub-model codes 48 or 54 associated with *-fose* in the ending table. In this case, there are four outcomes (codes 139, 141, 133, 135).

Regarding the other two possibilities, namely *F(-o)* and *Fos(-e)*, on comparing the sub-models codes of the stems *F-* and *Fos-* with those for the endings *-o*

and *-e*, no match is found. Therefore, these are not valid verbal forms.

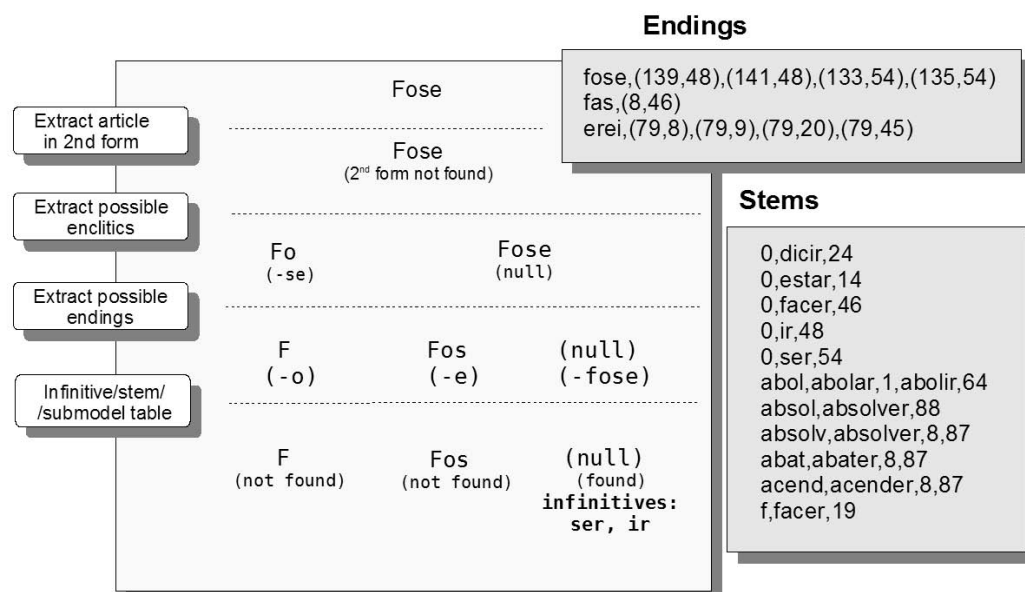


Fig. 7. Example of verbal analysis illustrating the different stages of this process (see text for details).

Figure 8 describes a more complex example where the verbal form includes a longer sequence of enclitics. In addition, an acute accent indicates a stressed syllable. The stress is maintained in the same syllable throughout the whole process, but the presence of the acute accent depends on acute accent rules for Galician. As can be observed in Figure 8, once the possible enclitics have been removed, we obtain four potential verbal forms that give rise to five stem-ending pairs. Although in two of these pairs the ending is null (no ending in the ending table matches the termination of the form under analysis), we cannot conclude that these verbal forms are not valid, because some complete forms of irregular verbs are directly included in the stem table. In this example, the final decision is that the analysed word is a compound of the verbal form *Fixera* and the sequence of two joined enclitics *-llelo*, after having checked the enclitic sequence is compatible with the resulting verbal form (as it may depend on the tense, mood, person, etc.).

Figure 9 shows an example containing the second form of the article, and so — after removing the second form — we rebuild the verbal form with a final *-s* or *-r*. Once the verb has been rebuilt the process is similar to the previous example, although it is sometimes necessary to modify the acute accent.



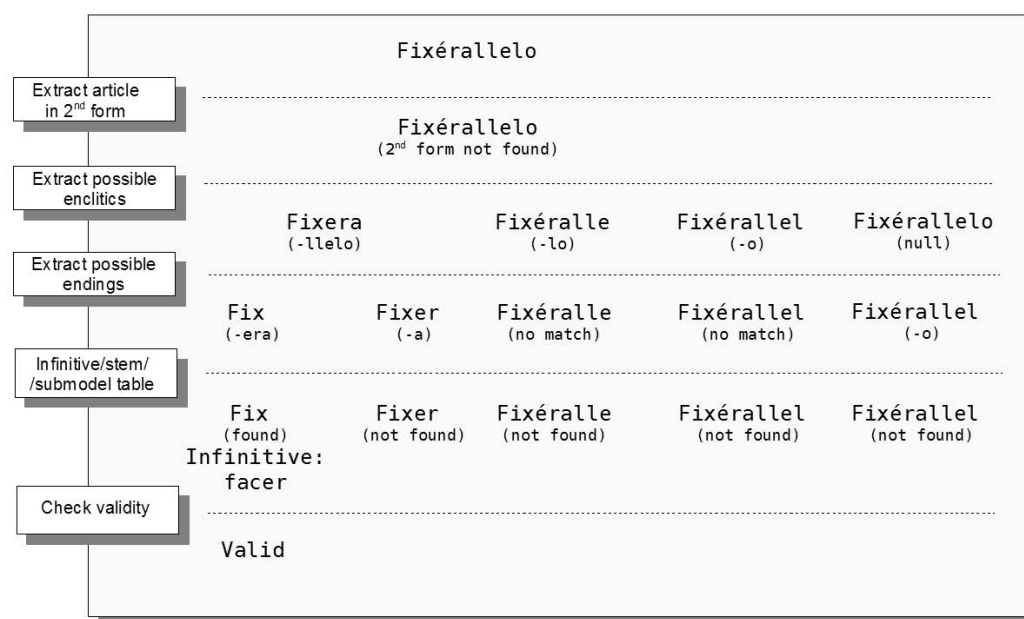


Fig. 8. Example of analysis of a verbal form with two enclitic pronouns (see text for details).

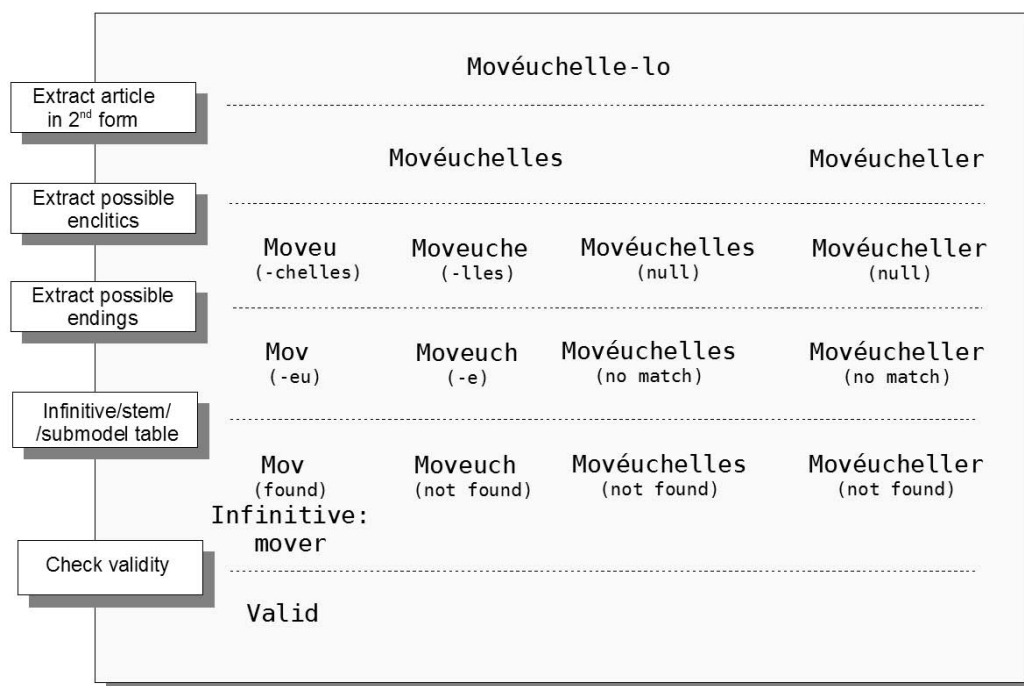


Fig. 9. Example of analysis of a verbal form with two enclitic pronouns and an article in second form (see text for details).

#### 4 Morphosyntactic analysis

Morphosyntactic analysis, or part-of-speech (POS) tagging, refers to the function that a word accomplishes within a sentence. The fact that most words

can fulfil several functions introduces ambiguity into tagging, and this problem is particularly aggravated in the case of Galician, given its wide use of contractions and clitics. Therefore, a strategy is required that will discern the right POS label for a word in a sentence that takes into account the influence of context.

There are two main approaches to POS tagging: using a large set of linguistic rules manually generated by a group of experts (Voutilainen et al. (1992)), or using data-based techniques which obtain the necessary information automatically from a previously labelled corpus (Brill (1995), Merialdo (1994)). Although rule-based taggers may be very accurate, developmental cost is generally very high, and this has led to most researchers focusing on the data-based techniques. In our case, we decided to follow a hybrid approach (Campillo (2005)) that involved allocating the tagging process to two stages: first, a small set of linguistic rules is used to reduce ambiguities (involving elimination of combinations that are clearly impossible), and second, a statistical model identifies the most probable sequence of categories in a sentence.

As mentioned previously, Galician is a minority language, which implies a serious lack of linguistic resources, including text corpora. It was therefore necessary to create a Galician corpus for POS tagging (Seijo et al. (2004)). Our text corpus was inspired by other corpora widely recognised as reliable for other languages, such as the Penn Treebank or Brown Corpus for English, Negra for German, and LexEsp for Spanish. It currently has about 450,000 words drawn from the journalistic corpus used in previous sections, to which no restrictions were applied in terms of style, topic or author. A bootstrapping technique was used to tag the texts comprising the corpus. First, a small amount of text was automatically tagged using the set of linguistic rules mentioned above and the result was manually revised by a team of linguists. The revised tagged text was then used to train the statistical model. The process continued with successive stages of tagging new text using the model, manually revising the outcome, and updating the statistical model. Since the statistical model accumulated information with each update, the number of errors to be manually corrected became successively smaller.

The design of the tagset used in the analysis was fundamental to satisfactory performance of the model, as this affects not only the utility of the results, but also the amount of tagged text needed to train the corpus. A small tagset would imply a small text corpus to obtain reliable statistics, but the analysis would provide less information. For our system we decided to use two different tagsets. The first one, designed by expert linguists, includes complete information for each word. It is composed of twelve main categories (noun, pronoun, adjective, determiner, article, verb, adverb, preposition, conjunction, interjection, punctuation and residual), plus the corresponding subcategories

(different types of adverbs, conjunctions, etc.). The full tagset includes about 130 categories, to which information about inflection, genre and number has been added. The second tagset — hidden to the user, since it is used only internally by the tagger — is a reduced version of the first tagset. It is composed of 49 categories (genre and number information included), obtained by clustering the categories of the main tagset according to common characteristics. Table 8 shows the mapping between the reduced and main tagsets. Regarding the design of the reduced tagset, worthy of comment is the handling of contraction categories, given their significance in Galician as opposed to other languages such as Spanish. Different experiments were carried out to check whether specific contraction categories should be added to the reduced tagset, but the best results were obtained when they were introduced, split into their different components, in the statistical model. This preserves the grammatical structure of sentences without the need to add new categories to the reduced tagset.

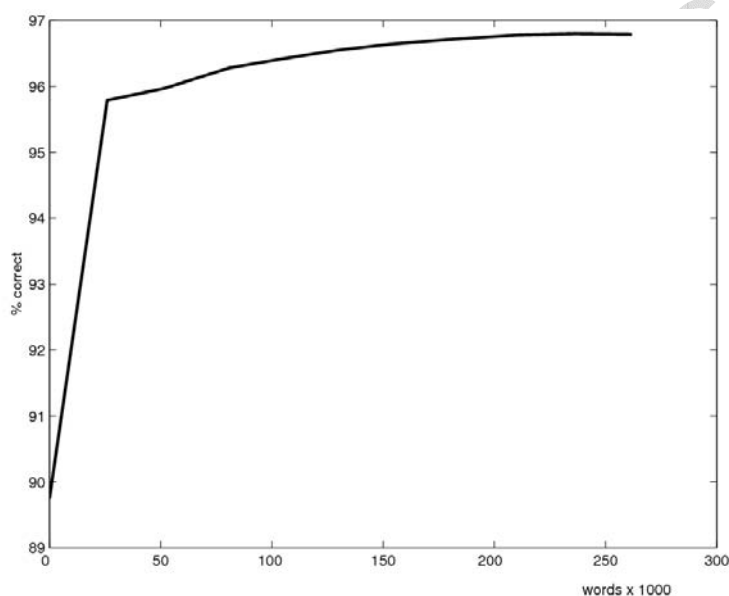


Fig. 10. Accuracy of the POS tagger as a function of the size of the training corpus (about 96% for only 50,000 words)

Figure 10 summarises accuracy data for the POS tagger as a function of training corpus size when used on an independent test text of about 25,000 words. As can be observed, accuracy is about 96% for only 50,000 words, which makes the tagger comparable to many other taggers. Table 9 shows the contribution of the different parts of the model. Firstly, demonstrated is the accuracy of the set of linguistic rules overriding the statistical model (89.58%), which was obtained by forcing the system to choose one of the categories not discarded by the rules. In our tagger, for implementation reasons that favour certain more likely categories, always choosing the first category obtains a slight improvement in accuracy (3.06%) with respect to random selection (86.52%). Secondly, demonstrated is the accuracy of the statistical model overriding the

Reduced tagset	Main tagset
Adjective	Adjectives and past participles.
Adverb	Adverbs and adverbial phrases
Exclamation mark (open)	“ ¡”
Question mark (open)	“ ¿”
Article	Indefinite and definite articles
Comma	“ ,”
Coordinate conjunction	Coordinate conjunctions
Subordinate conjunction	Subordinate conjunctions
Determiner	Demonstrative, indefinite and numeral determiners. Contracted indefinite.
Possessive determiner	Possessive determiners
Exclamative	Exclamative pronouns and determiners
Interrogative	Interrogative pronouns and determiners
Interjection	Interjections
Noun	Nouns and infinitives
Exclamation mark (close)	“ !”
Question mark (close)	“ ?”
Preposition	Prepositions and prepositional phrases
Pronoun	Stressed personal pronouns, indefinite, demonstrative and numeral pronouns. Contracted demonstrative and indefinite.
Possessive pronoun	Possessive pronouns
Unstressed pronoun	Unstressed pronouns
Period	“ .” , “ ;” , “ .”
Relative	Relatives
Verb	Verbs and verb phrases

Table 8

Mapping between the main and the reduced sets of labels.

linguistic rules (96.47%) which is clearly higher than that obtained with the linguistic rules alone. Finally, Table 9 shows the accuracy of the hybrid model (97.23%), which applies the statistical model to the categories not discarded by the linguistic rules. In this last case, the use of the linguistic rules results in a slight improvement that is explained by the ambiguity reduction in the

the statistical model inputs. Moreover, the linguistic rules prevent the tagger from occasional serious errors (from a strictly linguistic point of view).

Finally, if the full tagset (instead of the reduced set) were used in the statistical stage, the label accuracy of the hybrid model would drop to 93.38%. Thus, this result reinforces the decision to use an intermediate reduced tagset.

	Linguistic rules			Statistical model			Hybrid model		
	Label	Genre	Number	Label	Genre	Number	Label	Genre	Number
Test	89.58	86.46	87.98	96.47	97.69	98.56	97.23	97.93	98.78

Table 9

Accuracy comparison among the rule-based model, the statistical model and the hybrid scheme. Combining the rule-based and the statistical models into a hybrid scheme results in a significant improvement over the two models in isolation.

## 5 Intonation issues

Although it is beyond the scope of this article to describe Galician intonation, a brief comment on certain aspects may be of interest to the reader.

Unlike other languages such as English, where substantial research effort has been devoted to the phonological description of intonation and its application to text-to-speech systems (Pierrehumbert (1980)) no phonological study exists for Galician that is ample enough to be used in this type of system. Galician intonation shows some clear differences from Spanish. For instance, pitch range is usually greater in Galician than in Spanish. Moreover, in Galician the final part of the pitch contour in absolute interrogatives is descendant, in clear contrast with the final ascending pitch in this type of sentence in Spanish.

Due to the nonexistence of an exhaustive phonological description of Galician, in our *Cotovía* text-to-speech system we directly use a phonetic model of intonation. It is a classic hierarchical model in which the basic intonation unit is the accent group, considered as a sequence of unaccented words ending in an accented word. Intonation groups are formed by concatenating accent groups and sentences are composed of one or more intonation groups. This kind of intonation model has been previously proposed for other languages, among them Spanish (López (1993), Garrido (1996)). In our system, this intonation model has been successfully combined with unit selection techniques with the goal of providing our text-to-speech system with richer and more natural intonation (Campillo and Banga (2006)).

## 6 Conclusions

In this article we described some important aspects of the Galician language that need to be taken into account in speech technology development. In particular, we described language-specific problems encountered and the solutions developed for our text-to-speech synthesis system, especially those related to phonetic transcription and morphosyntactic analysis. This discussion was supported by numerous statistics, including relative frequencies for phonemes and diphthongs as well as the accuracy of the morphosyntactic analysis.

Certain characteristics of Galician have been excluded here because their treatment is similar to other languages. For instance — although there are some irregularities — the formation and detection of plurals can be basically implemented by means of a set of rules based on certain suffixes. The techniques described have been implemented in our text-to-speech system, *Cotovía*, and some of them are now being used for automatic translation research. The treatment of Galician vocalic elisions and assimilations has also been considered for alternative transcriptions in continuous speech recognition.

## 7 Acknowledgements

This work has been partially supported by the Spanish government, ERDF funds and the *Xunta de Galicia* under the projects TEC2006-13694-C03-03, HUM2005-08282-C02-01 and PGI0IT05TIC32202PR. *Cotovía* was developed in collaboration with the *Centro Ramón Piñeiro para a Investigación en Humanidades*. The authors wish to thank the anonymous reviewers for their valuable suggestions for improving the initial version of this paper.

## References

- Brill, E. (1995). Transformation-based error-driven learning and natural language processing: a case study in Part of Speech Tagging. *Computational Linguistics*, 21(4):543–565.
- Campillo, F. (2005). *Síntesis de voz basada en selección de unidades acústicas y prosódicas*. PhD thesis, ETSI Telecomunicación. Universidad de Vigo. Spain. Available from [http://www.gts.tsc.uvigo.es/web/imaxes\\_user/071023111001\\_tesis.pdf](http://www.gts.tsc.uvigo.es/web/imaxes_user/071023111001_tesis.pdf).
- Campillo, F. and Banga, E. (2006). A method for combining intonation modelling and speech unit-selection in corpus-based speech synthesis systems. *Speech Communication*, 48:941–956.

- Castro, O. (1998). *Aproximación a la Fonología y Morfología Gallegas*. University Press of the South, Inc.
- Garrido, J. M. (1996). *Modelling Spanish intonation for text-to-speech applications*. PhD thesis, Facultad de Lletres, Universitat de Barcelona, Spain.
- González, M. and Regueira, X. (1994). Estudio acústico das vocais galegas. In *Actas do XIX Congreso Internacional de Lingüística e Filoloxía Románicas*, volume 6, pages 141–179, A Coruña: Fundación Pedro Barrié de la Maza, Conde de Fenosa.
- González, M., García Mateo, C., Rodríguez Banga, E., and Fernández Rei, E. (2002). *Diccionario de verbos galegos Laverca*. Edicións Xerais de Galicia S.A.
- López, E. (1993). *Estudio de técnicas de procesado lingüístico y acústico para sistemas de conversión texto voz en Español basados en concatenación de unidades*. PhD thesis, E.T.S.I. de Telecomunicación, Universidad Politécnica de Madrid, Spain.
- Merialdo, B. (1994). Tagging English text with a probabilistic model. *Computational Linguistics*, 20(2):155–171.
- Moreno, A., Torre, D., Curto, N., and de la Torre, R. (2006). Inventario de frecuencias fonémicas y silábicas del castellano espontáneo y escrito. In *IV Jornadas en Tecnologías del Habla*, pages 77–81.
- Olive, J., Greenwood, A., and Coleman, J. (1993). *Acoustics of American English Speech*. Springer-Verlag.
- Pierrehumbert, J. (1980). *The Phonology and Phonetics of English Intonation*. PhD thesis, MIT, Cambridge, MA.
- Seijo, L., Martínez, A., Méndez, F., Campillo, F., and Banga, E. (2004). A Galician textual Corpus for morphosyntactic tagging with application to text-to-speech synthesis. In *Proceedings of LREC*, volume 5, pages 1759–1762, Lisboa.
- Veiga, A. (1976). *Fonología gallega*. Editorial Bello.
- Voutilainen, A., Heikilla, J., and Anttila, A. (1992). *Constraint Grammar of English. A Performance-Oriented Introduction*. Publication No 21, Department of General Linguistics, University of Helsinki.