# Phonetically Transcribed Speech Corpus Designed for Context Based European Portuguese TTS

Maria Barros[1] and Bernd Möbius[2]

[1] Eurescom GmbH, Wieblinger Weg 19/4,
69123 Heidelberg, Germany
[2] University of Bonn, Am Hof 1,
53113 Bonn, Germany
barros@eurescom.eu, moebius@ifk.uni-bonn.de

**Abstract.** This paper presents a speech corpus for European Portuguese (EP), designed for context based text-to-speech (TTS) synthesis systems. The speech corpus is intended for small footprint engines and is composed by one sentence dedicated to each sequence of two phonemes of the language, incorporating as many language contexts as possible at diphone and word levels. The speech corpus is presented in three forms: its orthographic sentences, its phonetic transcriptions considering the words coarticulation, typical phenomenon in EP and its phonetic transcriptions considering the words coarticulation and the vocalic reduction effects, which is another typical phenomenon of EP language. The paper describes the corpus and presents the results in total number of phonemes and their distribution in the word and the sentence, and the total number of diphones, considering and not considering the vocalic reduction for the phonetic transcription.

**Keywords:** Speech corpus, European Portuguese corpus, text-to-speech corpus, context based corpus.

## 1 Introduction

In a TTS system, the speech synthesis quality is strongly affected by the quality of the speech corpus, the quality of the natural language processing module and the quality of the speech generation module. The speech corpus quality, on its turn, depends on the quality of the corpus design, the quality of the annotations and transcriptions and the quality of the recording system.

Nowadays TTS systems are usually corpus based statistical systems. This means that they rely on the quality of the corpus used to produce better or worse results. Besides that, they are often systems using the language contexts for choosing the best speech units for a particular target. This makes it extremely important that the speech corpus design includes as many contexts as possible. Usually, statistical methods are used to search in large amounts of texts, in the intended speech units in the considered language contexts. For this reason the speech corpus for generic vocabulary systems are very large, sometimes too large to be considered.

The contextual factors that can be used are many. For example, the contextual factors used for the EP HMM-based speech synthesis system (Barros et al, 2005), considering the hierarchical structure, i.e., phone, syllable, word, phrase and utterance, were the following:

At phone level: current, previous and next phones; phones before previous phone and after next phone; and positions (forward and backward) of current phone in current syllable;

At syllable level: stress condition of current, previous and next syllables; number of phones in previous, current and next syllable; positions (forward and backward) of current syllable in current word and in current phrase; number of stressed syllables before and after current syllable in the current phrase; syllable counts between the previous stressed syllable and the current syllable and between the current syllable and the next stressed syllable in the utterance; and vowel of the syllable;

At word level: number of syllables in the current, previous and next words; and positions (forward and backward) of current word in current phrase;

At phrase level: number of syllables and number of words in current, previous and next phrases; and positions (forward and backward) of current phrase in the utterance;

At utterance level: number of syllables, words and phrases in the utterance.

As a solution towards the best quality language context based systems with small footprint engines, efforts were concentrated in a new, manually designed, speech corpus. The speech corpus intends to achieve a large amount of important language contexts within a limited size. To do so, one sentence was constructed targeting each sequence of two phonemes of the language. In this sentence, it was tried to include as many contexts as possible. Through statistical searches it is possible to obtain as many features as it is needed, but it would need a very large corpus and that is difficult for languages without many resources.

This speech corpus is designed for European Portuguese, although the idea can be extrapolated to other languages. The same data should be used to train the natural language processing (NLP) tasks and the synthesis engine, as this way achieves a better synchronization between the units selected by the NLP and the units that can be found in the synthesizer database.

## 2   Methodology

The methodology used to design the corpus took into consideration the 38 phonemes from EP, presented in figure 1, and the silence unit, to construct one sentence dedicated to each of these 39 units combined with each other and themselves. Excluding the combination of silence with itself there were 1520 possible sentences to be constructed. The result is a speech corpus with 1436 sentences, once the rest of the combinations were not possible due to language rules that will be explained in the description section.

The sentences were all manually constructed, trying to take into account a number of language context factors and considering, but not limiting to, the most used words in the lists relating to European Portuguese usage of the words, collected from the first Morpholympics for Portuguese (Linguateca, 2003). These lists were extracted from a set of 613 different texts, from different fields (literature, politics, advertising,

informal chats, general news, etc.), collected through different sources (newspapers, books, net, advertising, etc.), containing 80.903 text units correspondent to 17.128 different units. More about the description of the texts and the lists, with statistics, is available in the site (Linguateca, 2003).

| SAMPA | IPA |
|---|---|
| Oral Vowels and semi-vowels | |
| 6, a, E, e, @, O, o, u, i, j, w | ɐ, a, ɛ, e, i, ɔ, o, u, i, j, w |
| Nasal Vowels and semi-vowels | |
| 6˜, e˜, o˜, u˜, i˜, j˜, w˜ | ɐ̃, ẽ, õ, ũ, ĩ, j̃, w̃ |
| Fricative Consonants | |
| v, f, z, s, S, Z | v, f, z, s, ʃ, ʒ |
| Liquid Consonants | |
| L, l, l˥ | ʎ, l, ɫ |
| Vibrant Consonants | |
| r, R | r, R |
| Plosive Consonants | |
| b, p, t, k, g, d | b, p, t, k, g, d |
| Nasal Consonants | |
| m, n, J | m, n, ɲ |

**Fig. 1.** EP Phoneme Inventory, in SAMPA and in IPA

The contextual factors that were considered when building the sentences are:

At the diphone level, occurrences of the diphone: positioned at the beginning/end of the sentence; positioned at the beginning/middle/end of the word; and between two words.

At the word level, occurrences of the: words containing the target diphone at the beginning/middle/end of the sentence; and combinations with the target diphone between two words positioned at the beginning/middle/end of the sentence.

The speech corpus is presented in three forms:

- ➢ The orthographic sentences;
- ➢ A phonetic transcription considering the words coarticulation (that is a natural effect in EP continuous speech), following the grapheme-to-phoneme conversion (G2P) rules for the EP language;
- ➢ A phonetic transcription considering the words coarticulation and the vocalic reduction effects, common in EP language.



**Fig. 2.** Example of the sentences in the Speech Corpus

Figure 2 shows two examples of sentences, for the phonetic sequences /vu/ and /dg/, marking with circles all the possibilities of the sequence occurrence, independently of have been verified or not. In the example it is possible to see the difference in the number of the sequence occurrences for both types of phonetic transcription and verify that some diphones don't exist in the language but can appear when considering vocalic reduction.

## 3   Description

The speech corpus presented here has a total of 1436 sentences, comprising 5853 different words, with a total of 21500 words occurrences.

Two phonetic transcriptions for the orthographic sentences are provided, one following the G2P rules for EP language and the other considering the vocalic reduction effect that is common in the EP language. The number of words in the phonetic speech corpus transcribed following the G2P rules is 6264 and in the one taking the vocalic reduction into consideration, it is 6415. The difference in the number of words is due to the different transcriptions for some words when considering the vocalic reduction effect.

The effect of coarticulation between words is present in EP continuous speech. Regarding to this effect the same word can have different transcriptions, because there are graphemes that have different phonetic transcriptions according to the following one. For example, the grapheme <s> is transcribed as a /S/, if at the end of a sentence or followed by an unvoiced consonant (/p, t, k, s, S, f/ for EP); as a /Z/, if followed by a voiced consonant (/b, d, g, m, n, J, z, v, Z, l, l~, L, r, R/ for EP); or as a /z/, if followed by a vowel. Another example is the grapheme <l>, which is transcribed as a /l~/, if at the end of a sentence or followed by a consonant; or as a /l/, if followed by a vowel. One example of these situations in EP words is the word <vais>, which is transcribed as /vajS/, if it is followed by a word starting with an unvoiced consonant or if it is at the end of a sentence; as /vajZ/, if it is followed by a voiced consonant; and as /vajz/, if it is followed by a vowel. Another example is the word <mal>, which is transcribed as /mal~/, if it is followed by a word starting with a consonant or if it is at the end of a sentence; and as /mal/, if it is followed by a vowel.

The vocalic reduction is another effect present in EP continuous speech. It can be reflected by the suppression of /@/, by the phoneme /u/ reduction or suppression, or by the phonemes /u~/, /i/ and /i~/ reduction. Considering this effect, some language contextual factors are present that otherwise would not be found:

The phoneme /@/ suppression can happen in the middle of the words or at the end of the words or sentences. Some examples of this phenomenon are: <de> that would be transcribed as /d/ instead of /d@/; <dedica> that would be transcribed as /ddik6/ instead of /d@dik6/; <pote> that would be transcribed as /pOt/ instead of /pOt@/; and <apetece> that would be transcribed as /6ptEs/ instead of /6p@tEs@/.

The phoneme /u/ reduction or suppression happens in the middle of the words or at the end of the words or sentences. When the vowel /u/ suffers reduction it gives place to the semi-vowel /w/. An example of a word which usually presents /u/ suppression and reduction is the word <séculos>. Following the G2P rules this word is transcribed

as /sEkuluS/, but is commonly found as /sEklwS/ due to vocalic reduction. The paper (Barros et al, 2001) covers this phenomenon.

The phonemes /u~/, /i/ and /i~/ can suffer reduction, giving place to the semi-vowels /w~/, /j/ and /j~/, respectively, but they cannot suffer suppression. Examples of these situations are: <conjuntura>, which would be transcribed as /ko~Zw~tur6/ instead of /ko~Zu~tur6/, <litoral>, which would be transcribed as /ljtwral~/ instead of /litural~/, and <inventa>, which would be transcribed as /j~ve~t6/ instead of /i~ve~t6/.

In EP many of the consonants can not be found next to each other in the same word. Considering the case of the phoneme /@/ vocalic reduction means to include almost all the combinations between consonants that in other ways would not exist.
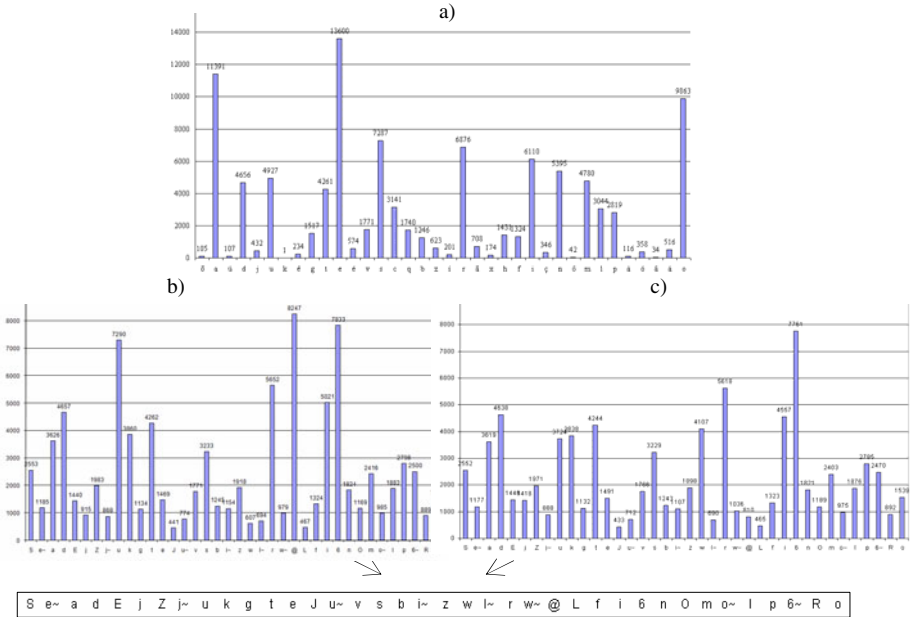


**Fig. 3.** Number of:  a) Grapheme occurrences; b) Phoneme occurrences, by rules; c) Phoneme occurrences, with vocalic reduction

### 3.1   The Speech Corpus by Graphemes

The 1436 sentences speech corpus comprises of 36 graphemes (considering the graphemes with accents - à, á, â, ã, é, ê, í, ó, ô, õ, ú - as individual units) and a total of 101750 graphemes occurrences. The total number of occurrences of each grapheme is shown in figure 3 a).

Although <k> is not a Portuguese grapheme it can appear in some foreigner words adopted in Portuguese vocabulary. In this case it appears in the word <snack>.

The <à> is a particular case in Portuguese language because it only appears with the contraction between <a a> and it's not possible in the middle of a word. For

example, the sentence <Dou à Joana.> (meaning <I give to Joana.>) has a grammar construction that comes from <Dou a a Joana.> (that would mean: <I give to the Joana.>). The examples present in the speech corpus are found in the words: <à>, <àquela>, <àquelas>, <àquele> and <àqueles>.

The <ç> is the only consonant that cannot appear in the beginning of a Portuguese word. This grapheme is always phonetically transcribed as /s/, which in the beginning of a word is produced by <s>, or <c> if followed by <(e, i)>.

The <ô>, <õ> or <ã> are the vowels that cannot appear in the beginning of a word. Relating to the last two, the nasal vowels, the <~> is used to nasalize the vowel and in cases with the vowel in the beginning of the word, the vowel is nasalized by using the consonants <(n,m)> in front of it.

The <o>, <a> and <s> are the most common endings for masculine words, feminine words and plurals, respectively. The infinitive form of Portuguese verbs always ends with <r>, what justifies the large amount of this grapheme's occurrences in the end of words and sentences.

From the Portuguese grammar rules, besides the consonants <s> and <r>, already mentioned above, only the consonants <m>, <z>, <l> and <b> can appear at the end of a word.

There are no cases of words ending with <â> or <õ>. In the last case, in EP the nasal sound in the end of the word is achieved by following the vowel with <m>.

The <h> is a special consonant in EP, once by itself does it not have a sonorous transcription, as it lost its aspiration with the evolution of the language. It is followed only by vowels, not by consonants, and it is used for the combinations <ch>, <nh> and <lh> that produce the phonemes /S/, /J/ and /L/, respectively.

## 3.2  Phonetic Transcription by Rules

The first phonetic transcription presented follows the G2P rules for the Portuguese language taking into consideration the words coarticulation natural in EP continuous speech, but not the vocalic reduction effect.

The 1436 sentences speech corpus comprises 38 phonemes and 92618 phonemes occurrences.

The total number of occurrences of each phoneme is shown in figure 3 b).

In European Portuguese language it is not possible to have words starting with the phonemes /l~/, /r/ or /J/, even considering phonetic effects inherent to continuous speech, unless it is some particular word like words imported from other Portuguese language dialects. This is the case of the word <nhangue>, imported from Angola Portuguese dialect to represent the name of an African bird and a place in Angola.

When considering the phonetic transcription following the G2P rules it is not possible to have words starting with a semivowel, although this can happen when considering the vocalic reduction effect.

All the vowels except /e~/ can be found at the end of a word or sentence. In the consonants case, the /d/, /k/, /g/, /t/, /J/, /v/, /s/, /L/, /f/, /n/, /m/, /l/, /p/ and /R/ are never found at the end of a word or sentence, unless due to particular cases of foreign words imported to European Portuguese vocabulary like, for instance, <internet>, <snack>, <camping>, <homeless> and <slogan>.

The consonants /Z/, /z/ and /l/ cannot be found at the end of a sentence or of an isolated word, but they can be found at the end of a word in continuous speech due to the words' coarticulation effect. The consonants /S/, /l~/ and /r/ can be found at the end of a sentence or isolated word.

Figure 5 presents a table with the total number of occurrences of each diphone when considering the phonetic transcription by rules, where the columns present the first phoneme of each diphone and the rows present the second one.

Any non-nasal vowel can only be next to a nasal vowel, or vice versa, if there is a small silence between words, in cases of slow rhythm of speech, because due to the words coarticulation effect a word that starts with a nasal vowel nasalizes the vowel from the precedent word and vice versa.

Most of the diphones composed by two consonants are not present in European Portuguese language, but there are some cases that are common: /pr/, /pl/, /br/, /bl/, /fr/ and /fl/; /Z/ followed by a voiced consonant; /S/ followed by an unvoiced consonant; /ks/ that transcribes the <x> for some grammatical cases (Barros and Weiss, 2006); /ps/, /pn/, /pt/, /kt/, /kn/, /gn/, /tn/, /gm/ and /tm/; /bt/, /bS/, /bZ/, /bz/, /bs/, /bv/ and /bm/; /dZ/, /dv/, /dm/ and /dr/; and /l~S/.

The <l> followed by a vowel is always transcribed as /l/, even if between words, unless there is a small silence in between the words, caused for instance by a comma, when then it is transcribed by a /l~/. It is also not possible to have /l/ or /l~/ followed by /r/, not because of the construction, that exists, but because in this case the <r> is transcribed by /R/ if in the same word, and there are no words starting with a /r/. Another case related to <l> not possible in EP is having /l/ or /l~/ followed by /J/, because <lnh> is not an allowed construction and there are no words starting with a /J/. Double "l", <ll>, doesn't exist in European Portuguese words, what makes the diphones among /l~/, /l/ and /L/ in the same word not possible. The /l~/ and the /L/ can never be preceded by consonant or a nasal vowel in the same word, as this would implicate a construction of three following consonants in a way that is not possible in EP language. Having a nasal vowel before a /l~/ implicates <(m,n)> after the vowel to nasalize it. So, with these cases, as well as with the consonants cases, it would be needed to have a consonant before the <l> and another after, to turn it in /l~/. It is also not possible to find these cases between words because there are no words starting with /l~/.

It's not possible to have /J/ preceded by a consonant, as it would lead to another type of construction of three consonants that is not allowed in European Portuguese language, the <nh> that gives the /J/ and any other consonant before. The same applies to nasal vowels preceding a /J/, because to have a nasal vowel before a <nh> it would have to be constructed by following the vowel with a <m> or <n>.

It's not possible to have /r/ in the same word after another /r/ or /R/, because the double <r> is always read as /R/. It is also not possible to have the /r/ after a nasal vowel in the same word, because it's the <n>, or the <m>, after a vowel that makes the nasal sound, but after a consonant the <r> is always read as /R/. There is no possibility of having these cases between words because there are no words starting with /r/.

It's not possible to have /n/ or /m/ preceded by a vowel and followed by a consonant, because in these cases the <m> and the <n> are used to nasalize the precedent vowel. Following this rule, it's not possible to have /n/ or /m/ followed by a consonant, because it would need a word construction of three consonants in a way that is not allowed in European Portuguese. There is one exception to these statements, which is the case of <m> followed by <n> and preceded by <a,o>, for example in the word <amnistia> or any word using the Portuguese prefix <omni>, which are transcribed as /6mniSti6/ and /Omni/, respectively.

The phoneme /e~/ is not found in the end of a word or after /e~/, /6~/ and /o~/, unless by influence of some regional accents/dialects, for instance in <lêem>, <têm> and <voem>, which are transcribed as /le~6~j~/, /t6~j~6~j~/ and /vo~6~j~/, but could be found in some regional accents as /le~e~/, /t6~e~/ and /vo~e~/, respectively.

### 3.3 Phonetic Transcription with Vocalic Reduction

The other phonetic transcription implemented to present the speech corpus considers the vocalic reduction effect that is common in the EP language, besides the words coarticulation. Due to the vocalic reduction there are less phonemes occurrences in this transcription than in the phonetic transcription by rules. The 1436 sentences speech corpus comprises of 38 phonemes and 84846 phonemes occurrences.

The total number of occurrences of each phoneme is shown in figure 3 c).

As it was explained before, in EP language it is not possible to have words starting with the phonemes /l~/, /r/ or /J/, even considering phonetic effects inherent to continuous speech, unless it is some particular word like words imported from other Portuguese language dialects.

Following the G2P rules it is not possible to have words starting with a semivowel, but this can happen when considering the vocalic reduction effect.

Due to the vocalic reduction effect, the vowel /@/ can be suppressed. This makes it possible to find almost all the consonants combinations. Also the consonant /S/ can be found before a voiced consonant from the transcription of <che> or <xe> that becomes /S/ instead of /S@/.

In European Portuguese language, the diphones /Ou/, /Ej/, /@j/, /aE/, /@E/, /aa/, /6a/, /ea/, /Ea/, /oa/ and /Oa/ aren't found in the same word, but can be found between words, both considering or not considering the vocalic reduction effect.

Figure 6 presents a table with the total number of occurrences of each diphone when considering in the phonetic transcription the vocalic reduction effect, where the columns present the first phoneme of each diphone and the rows present the second one.

The /@/ is suppressed in most of the cases of continuous speech, if considering the vocalic reduction effect, in any context of word or sentence.
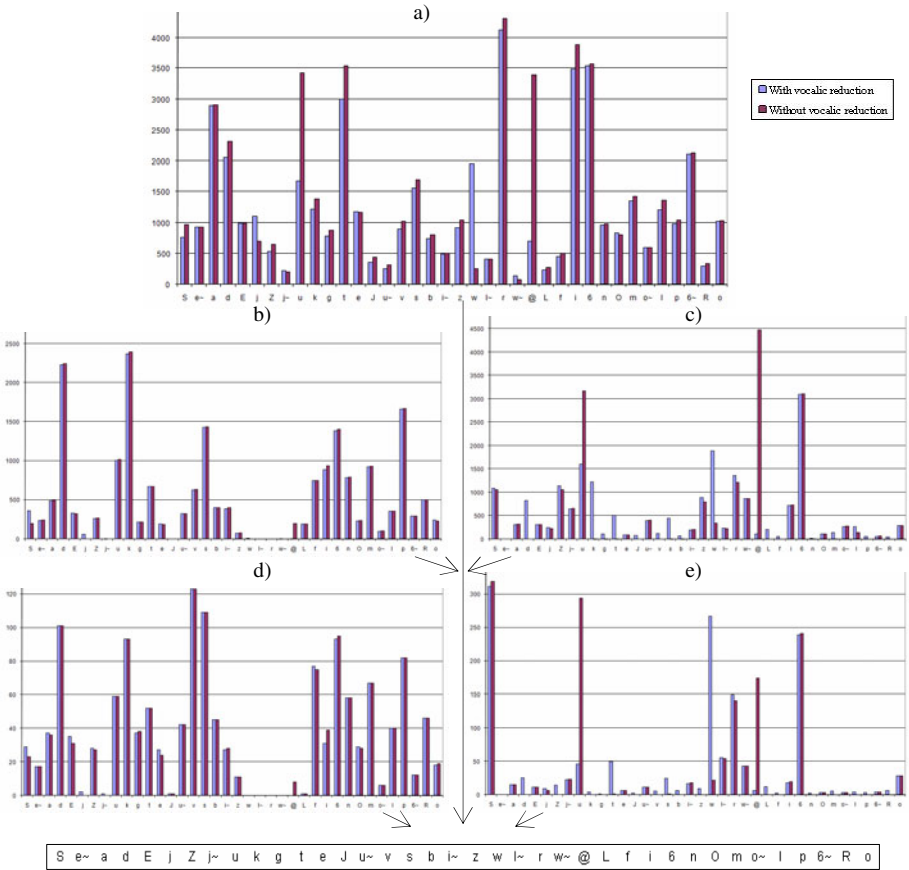
## 4   The Vocalic Reduction Influence

To better understand the influence of the vocalic reduction in EP continuous speech, a comparison between both phonetic transcriptions, following the G2P rules and

considering the vocalic reduction, in different contexts at word and sentence levels is presented.

Figure 4 a) presents a chart with the number of occurrences of each phoneme when in the middle of a word, for both phonetic transcriptions, following the G2P rules and considering the vocalic reduction. The same kind of chart is presented for the phonemes at the beginning of a word, in figure 4 b), at the end of a word, in figure 4 c), at the beginning of a sentence, in figure 4 d), and finally at the end of a sentence, in figure 4 e).



**Fig. 4.** Phonemes occurrences at the: a) Middle of word; b) Beginning of word; c) End of word; d) Beginning of sentence; e) End of sentence

From the figures it is possible to see that the biggest differences are in the number of occurrences of the phoneme /@/, the semi-vowels and the correspondent vowels. It is also possible to observe that the number of occurrences of some consonants in a particular context change if considering or not the vocalic reduction.

**Fig. 5.** Total number of diphone occurrences, by rules

Fig. 6 presents a 38 × 38 matrix of diphone occurrence counts. The row labels (first column) and column labels (header) use the same phoneme set. Best-effort transcription of the table follows.

| | o | R | 6ʲ | p | l | oʲ | m | O | n | 6 | i | f | L | @ | wʲ | r | lʲ | w | z | iʲ | b | s | v | uʲ | J | e | t | g | k | u | lʲ | Z | j | E | d | a | 6ʲ | S |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| S | 16 | 1 | 32 | 370 | 0 | 0 | 2 | 5 | 1 | 52 | 19 | 132 | 0 | 11 | 6 | 0 | 0 | 28 | 0 | 11 | 0 | 233 | 2 | 10 | 0 | 34 | 612 | 31 | 504 | 11 | 0 | 2 | 0 | 32 | 2 | 34 | 4 | 43 |
| 6ʲ | 0 | 0 | 34 | 98 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 7 | 0 | 0 | 0 | 1 | 0 | 0 | 6 | 0 | 23 | 112 | 18 | 1 | 0 | 0 | 678 | 9 | 45 | 0 | 0 | 2 | 0 | 0 | 85 | 88 | 2 | 9 |
| a | 5 | 8 | 2 | 49 | 161 | 3 | 46 | 6 | 17 | 15 | 13 | 47 | 36 | 0 | 0 | 812 | 333 | 238 | 152 | 9 | 79 | 89 | 149 | 11 | 1 | 0 | 105 | 45 | 34 | 8 | 0 | 15 | 230 | 6 | 606 | 10 | 68 | 53 |
| d | 126 | 45 | 64 | 99 | 32 | 10 | 75 | 37 | 39 | 818 | 431 | 65 | 4 | 67 | 2 | 28 | 0 | 758 | 74 | 44 | 31 | 82 | 31 | 27 | 3 | 125 | 76 | 16 | 122 | 381 | 2 | 97 | 64 | 68 | 83 | 287 | 4 | 89 |
| E | 5 | 21 | 3 | 30 | 106 | 4 | 36 | 14 | 25 | 23 | 15 | 39 | 8 | 2 | 0 | 267 | 75 | 8 | 14 | 16 | 58 | 122 | 64 | 23 | 0 | 6 | 103 | 85 | 64 | 37 | 0 | 64 | 1 | 6 | 48 | 14 | 32 | 59 |
| j | 7 | 32 | 5 | 22 | 25 | 2 | 14 | 16 | 15 | 177 | 15 | 14 | 16 | 9 | 2 | 131 | 0 | 36 | 72 | 12 | 6 | 30 | 3 | 4 | 1 | 5 | 70 | 5 | 34 | 95 | 0 | 24 | 16 | 4 | 36 | 113 | 31 | 287 |
| Z | 22 | 14 | 57 | 5 | 56 | 6 | 183 | 15 | 15 | 97 | 117 | 3 | 7 | 0 | 0 | 1 | 0 | 0 | 29 | 28 | 52 | 12 | 96 | 29 | 11 | 0 | 12 | 24 | 4 | 45 | 0 | 104 | 0 | 30 | 541 | 16 | 19 | 12 |
| u | 6 | 69 | 24 | 43 | 9 | 6 | 23 | 7 | 27 | 38 | 25 | 24 | 8 | 0 | 84 | 389 | 76 | 12 | 182 | 10 | 17 | 51 | 24 | 12 | 0 | 4 | 73 | 8 | 130 | 23 | 0 | 39 | 39 | 4 | 53 | 43 | 15 | 68 |
| k | 24 | 11 | 16 | 174 | 100 | 547 | 320 | 35 | 163 | 144 | 79 | 81 | 20 | 24 | 4 | 117 | 0 | 398 | 61 | 31 | 99 | 136 | 82 | 8 | 13 | 81 | 136 | 67 | 359 | 39 | 0 | 39 | 7 | 13 | 237 | 220 | 23 | 245 |
| g | 160 | 91 | 152 | 92 | 92 | 3 | 42 | 96 | 80 | 457 | 111 | 46 | 21 | 0 | 3 | 129 | 0 | 179 | 108 | 40 | 4 | 191 | 90 | 16 | 0 | 14 | 70 | 9 | 59 | 210 | 3 | 197 | 6 | 113 | 79 | 137 | 7 | 41 |
| t | 42 | 28 | 5 | 6 | 21 | 15 | 7 | 44 | 17 | 213 | 43 | 41 | 4 | 26 | 0 | 564 | 0 | 534 | 4 | 14 | 21 | 7 | 5 | 38 | 0 | 118 | 70 | 5 | 5 | 49 | 40 | 19 | 29 | 51 | 9 | 407 | 105 | 4 |
| e | 101 | 9 | 256 | 53 | 49 | 1 | 49 | 60 | 29 | 437 | 491 | 20 | 3 | 0 | 0 | 370 | 21 | 197 | 9 | 69 | 13 | 58 | 9 | 11 | 0 | 2 | 23 | 10 | 44 | 184 | 1 | 5 | 1 | 115 | 72 | 39 | 3 | 74 |
| J | 2 | 38 | 4 | 4 | 299 | 1 | 32 | 3 | 53 | 2 | 5 | 4 | 6 | 0 | 4 | 0 | 0 | 27 | 4 | 8 | 4 | 75 | 27 | 4 | 0 | 12 | 24 | 42 | 18 | 8 | 10 | 54 | 0 | 3 | 30 | 11 | 10 | 58 |
| uʲ | 22 | 8 | 40 | 3 | 10 | 5 | 2 | 10 | 2 | 132 | 11 | 28 | 3 | 0 | 3 | 31 | 0 | 0 | 12 | 4 | 39 | 14 | 7 | 24 | 0 | 142 | 93 | 2 | 4 | 13 | 0 | 68 | 1 | 5 | 4 | 170 | 63 | 2 |
| v | 7 | 3 | 8 | 37 | 9 | 3 | 23 | 48 | 9 | 17 | 14 | 5 | 3 | 55 | 0 | 12 | 0 | 77 | 3 | 17 | 5 | 55 | 6 | 7 | 12 | 164 | 6 | 17 | 32 | 36 | 0 | 5 | 59 | 8 | 98 | 182 | 90 | 27 |
| s | 94 | 11 | 106 | 8 | 33 | 92 | 7 | 114 | 13 | 219 | 305 | 12 | 2 | 64 | 0 | 233 | 53 | 184 | 19 | 52 | 15 | 17 | 24 | 5 | 0 | 56 | 22 | 3 | 6 | 129 | 0 | 22 | 178 | 133 | 13 | 10 | 24 | 19 |
| b | 156 | 8 | 353 | 44 | 46 | 68 | 35 | 25 | 46 | 262 | 361 | 54 | 4 | 0 | 2 | 1 | 53 | 93 | 15 | 139 | 24 | 53 | 48 | 16 | 0 | 6 | 13 | 63 | 75 | 71 | 9 | 14 | 8 | 126 | 57 | 140 | 73 | 9 |
| iʲ | 34 | 15 | 71 | 3 | 9 | 12 | 3 | 7 | 11 | 104 | 91 | 7 | 47 | 12 | 0 | 251 | 0 | 0 | 221 | 13 | 14 | 16 | 5 | 27 | 5 | 82 | 150 | 3 | 12 | 16 | 1 | 18 | 1 | 12 | 20 | 129 | 11 | 32 |
| z | 8 | 5 | 60 | 151 | 24 | 24 | 7 | 38 | 6 | 9 | 10 | 67 | 3 | 5 | 0 | 0 | 61 | 55 | 0 | 67 | 8 | 49 | 79 | 19 | 27 | 17 | 8 | 39 | 67 | 125 | 0 | 9 | 22 | 77 | 157 | 10 | 53 | 37 |
| w | 66 | 7 | 127 | 12 | 129 | 5 | 5 | 10 | 8 | 379 | 366 | 31 | 21 | 39 | 1 | 0 | 0 | 6 | 59 | 16 | 8 | 10 | 83 | 0 | 0 | 0 | 132 | 3 | 12 | 89 | 0 | 41 | 6 | 49 | 322 | 341 | 34 | 378 |
| lʲ | 28 | 8 | 34 | 282 | 4 | 0 | 122 | 77 | 173 | 148 | 145 | 62 | 6 | 0 | 0 | 290 | 0 | 248 | 5 | 97 | 33 | 158 | 30 | 35 | 0 | 83 | 281 | 4 | 295 | 0 | 0 | 17 | 46 | 89 | 71 | 27 | 5 | 8 |
| r | 0 | 76 | 5 | 46 | 21 | 20 | 178 | 4 | 124 | 0 | 0 | 28 | 9 | 0 | 0 | 1 | 0 | 0 | 23 | 20 | 21 | 202 | 61 | 12 | 29 | 3 | 36 | 128 | 62 | 417 | 103 | 237 | 4 | 20 | 242 | 34 | 6 | 16 |
| wʲ | 37 | 9 | 320 | 140 | 8 | 1 | 35 | 37 | 40 | 1029 | 479 | 31 | 2 | 2 | 1 | 85 | 0 | 64 | 2 | 2 | 24 | 65 | 4 | 47 | 0 | 0 | 19 | 97 | 352 | 51 | 1 | 80 | 0 | 14 | 130 | 32 | 12 | 77 |
| @ | 15 | 35 | 80 | 78 | 11 | 3 | 4 | 91 | 12 | 68 | 34 | 16 | 1 | 10 | 0 | 225 | 0 | 78 | 237 | 61 | 8 | 36 | 6 | 5 | 130 | 11 | 2 | 5 | 87 | 64 | 712 | 4 | 1 | 32 | 60 | 146 | 12 | 23 |
| L | 1 | 21 | 1 | 17 | 5 | 5 | 7 | 19 | 6 | 2 | 3 | 5 | 52 | 4 | 0 | 593 | 0 | 165 | 8 | 25 | 4 | 8 | 273 | 66 | 176 | 23 | 8 | 20 | 34 | 17 | 2 | 42 | 10 | 29 | 25 | 52 | 7 | 16 |
| f | 9 | 13 | 7 | 8 | 38 | 5 | 6 | 30 | 7 | 66 | 21 | 92 | 49 | 112 | 11 | 15 | 49 | 5 | 60 | 51 | 60 | 13 | 195 | 15 | 0 | 10 | 235 | 1 | 11 | 38 | 0 | 11 | 0 | 41 | 11 | 247 | 9 | 106 |
| i | 148 | 9 | 15 | 2 | 128 | 14 | 222 | 55 | 230 | 138 | 264 | 189 | 0 | 0 | 1 | 272 | 0 | 77 | 7 | 33 | 208 | 273 | 1 | 20 | 1 | 2 | 321 | 14 | 5 | 118 | 0 | 4 | 407 | 6 | 453 | 5 | 6 | 8 |
| 6 | 38 | 4 | 11 | 129 | 208 | 0 | 500 | 3 | 276 | 437 | 26 | 10 | 46 | 27 | 910 | 20 | 0 | 0 | 21 | 9 | 3 | 521 | 29 | 6 | 0 | 38 | 1 | 10 | 411 | 47 | 0 | 76 | 3 | 48 | 488 | 133 | 35 | 45 |
| n | 20 | 43 | 42 | 427 | 4 | 0 | 4 | 55 | 3 | 219 | 186 | 33 | 37 | 101 | 2 | 0 | 0 | 206 | 9 | 44 | 65 | 16 | 15 | 2 | 0 | 0 | 66 | 38 | 538 | 9 | 0 | 7 | 22 | 5 | 2 | 166 | 35 | 10 |
| O | 44 | 129 | 243 | 1 | 57 | 13 | 53 | 5 | 57 | 337 | 221 | 8 | 2 | 106 | 0 | 2 | 0 | 0 | 10 | 6 | 5 | 47 | 47 | 9 | 7 | 95 | 28 | 12 | 4 | 140 | 0 | 26 | 9 | 89 | 43 | 2 | 325 | 28 |
| m | 5 | 2 | 3 | 52 | 10 | 0 | 15 | 66 | 48 | 4 | 2 | 125 | 2 | 0 | 2 | 456 | 0 | 96 | 15 | 44 | 38 | 15 | 5 | 12 | 0 | 1 | 143 | 23 | 37 | 233 | 0 | 10 | 29 | 78 | 30 | 45 | 8 | 13 |
| oʲ | 32 | 33 | 104 | 12 | 4 | 2 | 11 | 94 | 6 | 538 | 176 | 7 | 1 | 89 | 0 | 0 | 0 | 220 | 3 | 22 | 10 | 109 | 9 | 0 | 0 | 54 | 9 | 4 | 43 | 310 | 0 | 32 | 9 | 6 | 27 | 10 | 15 | 39 |
| l | 6 | 5 | 5 | 80 | 2 | 26 | 25 | 6 | 2 | 2 | 3 | 4 | 3 | 2 | 0 | 301 | 0 | 8 | 8 | 10 | 30 | 38 | 6 | 2 | 15 | 187 | 24 | 0 | 17 | 9 | 0 | 9 | 44 | 25 | 167 | 0 | 8 | 29 |
| p | 63 | 7 | 60 | 14 | 4 | 38 | 5 | 16 | 9 | 88 | 70 | 4 | 16 | 0 | 0 | 0 | 0 | 30 | 42 | 12 | 10 | 80 | 48 | 7 | 0 | 3 | 341 | 0 | 17 | 34 | 0 | 23 | 149 | 5 | 17 | 0 | 0 | 0 |
| 6ʲ | 131 | 6 | 20 | 5 | 57 | 0 | 2 | 6 | 1 | 55 | 11 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 9 | 82 | 40 | 0 | 0 | 0 | 30 | 12 | 0 | 27 | 34 | 0 | 15 | 0 | 0 | 70 | 0 | 0 | 0 |
| R | 2 | 4 | 1 | 49 | 39 | 21 | 11 | 0 | 6 | 0 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 30 | 51 | 6 | 9 | 0 | 3 | 24 | 0 | 34 | 34 | 0 | 23 | 44 | 25 | 17 | 45 | 15 | 13 |
| o | 33 | 2 | 10 | 18 | 18 | 2 | 25 | 0 | 6 | 0 | 0 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 10 | 10 | 51 | 48 | 12 | 15 | 30 | 12 | 0 | 70 | 45 | 0 | 15 | 149 | 5 | 70 | 10 | 8 | 29 |

**Fig. 6.** Total number of diphone occurrences, considering the vocalic reduction

## 5   Conclusions

The paper presents a speech corpus for EP, for small footprint context based text-to-speech TTS systems, with a reduced size, but composed by sentences incorporating each diphone of the language in as many language contexts as possible. Besides considering diphones, it also considers sequences of phonemes that would not constitute a diphone of EP, but appear in EP continuous speech due to the vocalic reduction effect.

The speech corpus is presented in its orthographic sentences and in its phonetic transcriptions considering the words coarticulation only and also the vocalic reduction effects. An analysis between both phonetic transcriptions is presented, in order to understand the importance of this effect in EP. The contextual factors considered for the corpus design were the occurrences of each diphone positioned at the beginning and end of the sentence, or at the beginning, middle and end of the word, and between two words, at the diphone level, and the occurrences of the words containing the target diphone at the beginning/middle/end of the sentence, or combinations with the target diphone between two words positioned at the beginning/middle/end of the sentence, at the word level.

## References

1. Barros, M.J., Maia, R., Tokuda, K., Freitas, D., Resende, F.G.: HMM-based European Portuguese Speech Synthesis. In: Interspeech 2005, Lisbon, Portugal (2005)
2. Linguateca: The first evaluation contest for morphological analysers of Portuguese (2003), http://www.linguateca.pt/Morfolimpiadas/ (last visited on 02-09-2009)
3. Barros, M.J., et al: Backclose Nonsyllabic Vowel [u] Behavior in European Portuguese: Reduction or Supression. In: ICSP 2001, Taejon, South Korea (2001)
4. Barros, M.J., Weiss, C.: Maximum Entropy Motivated Grapheme-to-Phoneme, Stress and Syllable Boundary Prediction for Portuguese Text-to-Speech. In: IV Biennial Workshop on Speech Technology, Zaragoza, Spain (2006)