

Quantifying Temporal Speech Reduction in French using Forced Speech Alignment

Martine Adda-Decker and Natalie D. Snoeren

LIMSI-CNRS
B.P. 133, 91403 Orsay, France
`{madda,nsnoeren}@limsi.fr`

Abstract

The processing of speech reduction remains one of the major challenges to automatic speech recognition (ASR) systems, as speech reduction often results in a considerable number of automatic transcription errors. Conversely, automatic transcription errors may indicate interesting reduction phenomena. The present article focuses on temporal speech reduction in spoken French. In a series of explorations, we examined large speech corpora with respect to production variation in different speech styles, and in particular, to shorter pronunciations and disappearing sounds. An ASR tool was used, forced speech alignment, that allows one to locate and to quantify speech regions prone to temporal reductions. Our study made use of various styles of large speech corpora, including broadcast news, as well as telephone and face-to-face conversations, thereby including both casual and careful speech. The results highlight the increasing impact of temporal speech reduction with less formal, more spontaneous and interactive speaking styles. In a broader sense, our study provides a demonstration of how ASR systems can be employed to consistently explore variations in speech in virtually unlimited speech corpora.

Key words: Automatic speech recognition, forced alignment, speech reduction in French, pronunciation variants, speech style, segment duration distribution.

1. Introduction

2 An important bottleneck to a large-scale expansion of automatic speech
3 recognition (ASR) devices is the efficient processing of spontaneous or ca-

4 casual speech. Within the ASR community “spontaneous” generally refers to
5 all kinds of phenomena that make the speech signal deviate from a care-
6 fully articulated sequence of words and sounds. Problematic topics in casual
7 speech include word fragments and various sorts of disfluencies (Shriberg,
8 1994), speech sounds that overlap with other sounds, such as laughter or
9 crying, produced by the same speakers, phonological variants, underarticu-
10 lated speech and, last but certainly not least, speech reductions resulting in
11 missing sounds. Many speech scientists share the belief that much knowledge
12 can be gained from studying characteristics of casual speech (Greenberg &
13 Chang, 2000; Greenberg, Carvey, Hitchcock, & Chang, 2003; Nakamura, Fu-
14 rui, & Iwano 2006; Strik, Elffers, Bavcar, & Cucchiarini, 2006). For instance,
15 Greenberg and colleagues (2000, 2002) have investigated syllabic structures
16 in casual speech from the SwitchBoard data, Nakamura et al. (2006) com-
17 pared spectral properties of careful and casual speech on large Japanese cor-
18 pora, thereby highlighting spectral reduction. Strik and colleagues (2006)
19 have studied reduction phenomena in Dutch, with a focus on the problem of
20 disappearing sounds, especially in multiword expressions.

21 Speech reductions seem to first affect the least informative speech por-
22 tions (Jurafsky, Bell, Gregory & Raymond, 2001), i.e., function words that
23 are predictable from the context, idioms, morphological items (in particu-
24 lar endings), dates, discourse markers etc. Speech reduction produces either
25 different (centralised) phonemes, fewer phonemes, or even fewer syllables
26 (Ernestus, 2000; Van Son & Pols, 2003; Duez, 2003; Adda-Decker, Boula de
27 Mareuil, Adda, & Lamel, 2005).

28 As far as phonemic segmentation and labelling is concerned, it is far
29 from being obvious that an automatic speech recognizer will prefer the same
30 options as a human expert. A human listener can not always tell for sure
31 whether a phoneme is deleted since some of the missing phoneme’s acoustic
32 features may be present in adjacent phonemes, and may even be perceived.
33 Moreover, it is well-known that human speech perception may sometimes
34 be biased by higher-level language knowledge and understanding (see, e.g.,
35 Ganong, 1980; Elman & McClelland, 1988; Samuel & Pitt, 2003). A given
36 ASR system, on the other hand, will consistently make the same decisions
37 over the entire corpus, and can be parameterized to best fit the investigator’s
38 needs.

39 Many studies have addressed the issue of phone boundary reliability, as
40 well as agreement between several manual and/or automatic annotations.
41 In early work, Cole, Oshika, Noel, Lander and Fanty (1994) observed that

42 boundary location achieved 80% inter-annotator agreement regarding man-
43 ually labelled phone boundaries with a 10 ms tolerance. More recently, over
44 90% agreement with a 20 ms tolerance has been reported between several
45 automatic alignments and manual labelling with a 20 ms tolerance (cf. Ho-
46 som, 2009). The overall reliability of automatic labelling through forced
47 alignment can be considered close to the one achieved by human experts.
48 Following these results, ASR systems can be used for a large panel of em-
49 pirical studies in that they allow one to consistently investigate variations in
50 large speech corpora in terms of known influential parameters, such as speak-
51 ing style, gender, dialectal accents, and emotions. In this paper, however,
52 we propose a method to provide evidence of temporal speech reduction with
53 the help of global descriptors, such as phone segment duration distributions.
54 Increasing proportions of short segments are considered as indicative of a
55 higher density of temporal reduction. Corresponding speech regions need to
56 be further studied in order to gain deeper insight in the complexity of sponta-
57 neous speech specific reduction phenomena, to increase our understanding of
58 the general mechanisms underlying pronunciation variation and last but not
59 least, to contribute to better acoustic speech models for ASR in the future.

60 In the current paper, we are concerned with uncovering some of these
61 processes, in particular temporal speech reduction in French. The goal of the
62 present research is to identify speech regions that are prone to reduction using
63 the forced speech alignment tool (Adda-Decker & Lamel, 1999) based on the
64 LIMSI speech recognition system (Gauvain, Lamel, Adda, & Adda-Decker,
65 1994). It is demonstrated that forced speech alignment can be employed to
66 quantify speech reduction in large spoken corpora. We carried out the inves-
67 tigations with a special focus on various speech styles, ranging from broadcast
68 news to telephone and face-to-face conversations. By using large speech cor-
69 pora, we aimed to find out whether the extent to which speech reduction is
70 observed varies as a function of different speech styles and languages (French
71 and English). Moreover, we looked into the question of whether vowels and
72 consonants are more or less prone to temporal reduction. The final section
73 summarizes the main results and discusses the implications of the outcomes
74 of the corpus-based study. Before we turn to our corpus-based study, how-
75 ever, we will first provide a short overview of speech reduction and the type
76 of ASR errors it may give rise to in the French language, as the proposed
77 investigations are originally motivated by ASR error analyses.

78 2. Speech reduction and ASR errors

79 2.1. Speech reduction

80 A considerable amount of research has been devoted to the study of speech
81 reduction phenomena, including consonant lenition, consonant cluster sim-
82 plications, vowel reduction and syllable restructuring (see, e.g., Van Son &
83 Pols, 2003; Duez, 2003; Dilley & Pitt, 2007; Ernestus, 2000; Tseng, 2005).
84 Temporal structure reduction is frequently observed in spoken English (e.g.,
85 *isn't it* or *it's*) and spoken German (*ins* instead of *in das*, 'in the'). In French,
86 similar reduction phenomena occur. However, in the remainder of this paper
87 we are solely concerned with less explicit temporal reduction phenomena.
88 Such reduced pronunciations are generally not reflected in normative written
89 sources: *il y a* [ilija] ('there is') is most often uttered as *y a* [ja], and *je ne*
90 *sais pas*, [ʒənəsɛpa] ('I don't know') may have acoustic realisations close to
91 [ʃɛpa] or even [ʃpa], where the *ne* in the negative form *ne ... pas* is being
92 omitted and /ʒə/ and /s/ are merged to form a mere fricative segment with
93 a [ʃ]-like sound. Moreover, the /ɛ/-vowel may become devoiced and merged
94 with the preceding fricative segment. As these examples illustrate, the scope
95 of sequential reductions in French often surpasses word boundaries. Typi-
96 cally one or more short function words are involved. One common means of
97 addressing this problem from an ASR perspective in speech processing is by
98 adding "multiwords" in the pronunciation dictionary for observed sequences
99 of words that tend to co-occur more frequently than chance (see Strik &
100 Cucchiaroni, 1999; Strik, Binnenpoorte, & Cucchiaroni, 2005). For English,
101 *want to* can thus accept a pronunciation variant like [wʌnə] and for French
102 *je ne sais pas* can even receive a [ʃpa] reduced pronunciation variant. Strong
103 temporal reductions may also appear within groups of content words, such
104 as dates and compound nouns, as is illustrated by an example in English
105 (see Figure 1). The temporally reduced portion *student athletes* is further
106 detailed in Figure 2: a manual phonetic transcription of *student* is given
107 together with a canonical full form obtained by forced alignment. We will
108 come back to this example in the Method section. Before turning to our
109 corpus-based study and the related methodology, we will first discuss some
110 of the most frequent transcription errors that are observed for French.

111 2.2. Typical transcription errors in French

112 Previous studies have reported about 10% word error rates for French
113 careful (i.e., journalistic) speech and above 15% for casual telephone speech,

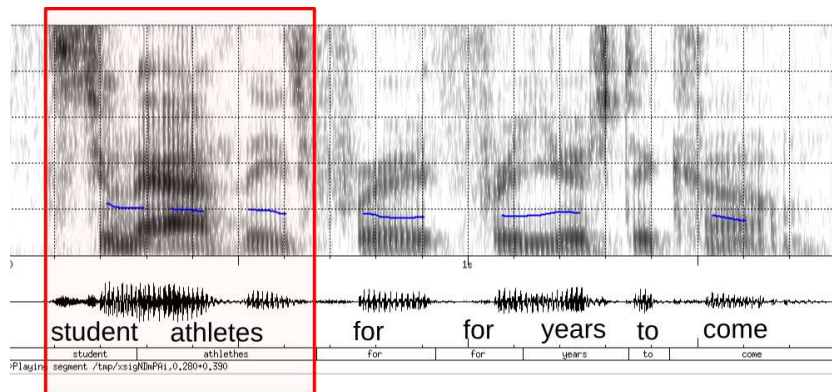


Figure 1: Speech signal of a reduction phenomenon in an English compound: *student athletes* /stjudənt æθlits/, is approximately being produced as [stjunæθlits] within the carrier sentence *it's gonna benefit student athletes (for) for years to come*.

Reference	Pronunciation	Hypothesis	Pronunciation
a	/a/	à	/a/
sait	/sɛ/	c'est	/sɛ/
cesse	/sɛs/	seize	/sɛz/
rentre	/ʁɑ̃trə/	rendre	/ʁɑ̃drə/

Table 1: (Near-)homophone substitutions in French that do not contain temporal reduction. Both the reference words and their transcription output hypotheses are indicated with their corresponding pronunciations.

114 with hundreds of hours of appropriate casual speech data and complex sys-
 115 tem combinations (see Lefèvre, Gauvain & Lamel 2005; Prasad et al, 2005).
 116 Among automatic transcription errors in careful speech, approximatively 30-
 117 40% of errors consist of homophone or near-homophone errors without tem-
 118 poral reduction. In Table 1, some examples are given of typical confusions
 119 that involve frequent and less frequent (near)-homophones. Reference words
 120 are substituted by the more frequently occurring homophone hypotheses.
 121 For example, *sait* ('knows') may thus be replaced by *c'est* ('that is'), the
 122 latter having a higher prior probability in the ASR language model. Table 2
 123 shows examples of near-homophone transcription errors due to typical tem-
 124 poral reduction phenomena in French, including cross-word vowel merging,
 125 word-final consonant cluster simplification and short function word deletion.

126

127 The above-mentioned examples correspond to journalistic, i.e. carefully

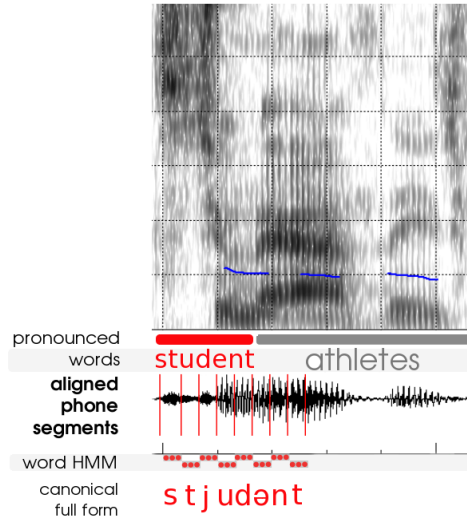


Figure 2: Zoom on the mismatched automatic alignment of the temporally reduced speech portion from the previous Figure. The noun *student* approximately being produced as [stjun] corresponds to the red bar portion. This is much shorter than the minimal duration of the automatically aligned word HMM based on the canonical full form /stjudənt/.

<i>Reference</i>	<i>Pronunciation</i>	<i>Hypothesis</i>	<i>Pronunciation</i>
ça avait	/saavɛ/	savaient	/savɛ/
semble que	/sɑ̃bləkə/	somme que	/sɔmkə/
parce que	/paʁsəkə/	ce que	/səkə/
près de Paris	/pʁɛdəpaʁi/	préparé	/pʁɛpaʁɛ/

Table 2: Near-homophone errors with temporal reduction: the system’s *Hypothesis* pronunciation is shorter than the *Reference* pronunciation.

128 prepared speech. When analysing casual, conversational speech, however, the
 129 proportion of errors due to temporal reduction increases significantly. The
 130 proportion of short word sequences, in particular discourse markers (*tu sais*,
 131 *tu vois* (‘you know’, ‘you see’)) and markers of reported speech (*il m’a dit*, *je*
 132 *lui ai dit* (‘he told me’, ‘I said to him’)) is particularly high. Consequently,
 133 these sequences are often prone to recognition errors, unless specific acoustic
 134 models have been trained on these specific sequences.

135 Finally, in French the schwa is known to contribute to temporal structure
 136 variation, giving rise to a wealth of automatic transcription errors. In Table 3,
 137 some illustrations of errors are listed that are due to insertions and deletions
 138 of the French schwa (see Verney Pleasants, 1956; Dausès, 1973; Fougeron,

139 Goldman & Frauenfelder, 2001; Adda-Decker, 2007). The listed examples
 140 are chosen from the French ESTER-2005 ASR system evaluation campaign
 141 and nicely illustrate the impact of relatively simple structure variations on
 142 recognition performance. About 5% of the errors can be related to appearing
 143 or disappearing schwas, leading to sequences that are simply incompatible
 144 with what was actually predicted by the acoustic word models.

Reference word string	Decoded word string	Comment
Absence of schwa in dict. & acoustic model, but produced by speaker		
<i>Marc <u>B</u>londel</i>	marque Blondel	consonant cluster
<i>en fait</i>	en fait de	final release
<i>le week-end pascal</i>	le week-end pascal le	phrase boundary
Presence of schwa in dict. & acoustic model, but not produced by speaker		
<i>tout <u>le</u> temps</i>	tout _ temps	idiom
<i>temps <u>de</u> leur installation</i>	temps _ leur installation	noun phrase
<i>quai <u>de</u> Seine</i>	quête saine	compound + assimilation
<i>c'était <u>le</u> même marasme</i>	c'est elle même marasme	homophone
<i>confiance appréciable <u>le</u> tandem</i>	confiance appréciable _ tandem	homophone

Table 3: Transcription errors arising from mismatches in temporal structures between the observations and the model due to word-final schwa. In the upper panel, the schwa was missing in the model, but observed in the speech signal. The lower panel illustrates the reverse situation.

145 The examples in the top panel show contexts where a schwa is being
 146 produced in the speech signal without any evidence in the written forms.
 147 The epenthetic schwa is due to contextual effects such as the creation of a
 148 complex consonant cluster at word boundaries (here /ʁkbl/ of the proper
 149 name *Marc Blondel*), or phrase-final releases. In this particular case, the
 150 French language permits near-homophone insertions of short and frequent
 151 function words such as *de* ('of') and *le* ('the').

152 The examples in the lower panel correspond to the more frequent situation
 153 of a mismatch due to shorter productions on the one hand and longer models
 154 on the other hand. An example worth mentioning is *quai de Seine* ('Seine
 155 bank') that illustrates the process of schwa deletion (of the word *de*) before
 156 a consonant (/s/ of *Seine*) after an open syllable (/kɛ/ of *quai*). Moreover
 157 the /d/ is assimilated to [t] due to the following unvoiced /s/ (cf. Snoeren,
 158 Hallé & Segui, 2006). In casual speech, the situation of complex combinations

159 of various reduction processes appear to be very common. Using large cor-
 160 pora therefore enables us to give a synthetic and exhaustive overview of the
 161 various reduction processes (cf. Schuppler, Ernestus, Scharenborg & Boves,
 162 2008). As a first step in this direction, we propose to quantify temporal re-
 163 ductions using forced alignment and canonical pronunciations. This allows
 164 us to measure deviations from canonical temporal structures in terms of their
 165 phone segment duration distributions.

166 3. A corpus-based study using forced alignment

167 Forced speech alignment consists in linking a reference transcription to
 168 its acoustic speech signal using an ASR system. The resulting word and
 169 subword (typically phoneme) boundaries are determined with respect to the
 170 ASR system’s configuration (acoustic phone models, model topology, word
 171 pronunciations, inter-word optional silences and so forth). Forced alignment
 172 is typically used to automatically label large manually transcribed speech
 173 corpora for acoustic model training. The resulting phone labels and bound-
 174 aries are not necessarily in line with manual phonetic segmentation, nor com-
 175 pletely compatible with different system configurations. However, previous
 176 studies have shown that major linguistic trends (e.g. vowel reduction and
 177 duration) can be consistently observed whilst using ASR systems developed
 178 independently by different research teams (Adda-Decker, Gendrot & Nguyen,
 179 2008).

180 In a first series of explorations, we compared different speaking styles
 181 using phone segment duration distributions as obtained by forced alignment
 182 (Adda-Decker & Lamel, 2005). The questions we were interested in are the
 183 following. First, what is the effect of casual speech on the duration distribu-
 184 tion as compared to careful speech? Second, how does the French data-set
 185 compare to the English data-set? Third, do the observed results hold for
 186 different types of casual speech? Finally, does the extent to which speech
 187 reduction occurs vary for vowels and for consonants? For this latter compar-
 188 ison between vowels and consonants, all segment duration distributions were
 189 examined using the following four duration classes:

	short:		\leq	40 ms
	medium:	50	-	110 ms
190	long:	120	-	240 ms
	very long:		\geq	250 ms

191 Similar duration classes have proven to be useful in showing the impact
 192 of duration on F1/F2 formant values of oral vowels (see Gendrot & Adda-
 193 Decker, 2005). The *very long* duration class was mainly introduced to check
 194 the proportion of very long segments and includes silence and sounds other
 195 than speech (e.g., hesitations). Most temporal reductions of interest are
 196 assumed to be found in the duration class labelled *short*. It is important
 197 to point out that the tuning of 40 ms as upper duration limit should be
 198 considered as a permissive one, and by no means as a norm. On the basis
 199 of our speech data, it has been observed that a tighter limit of 30 ms would
 200 reduce the proportion of segments in this class by 50% for careful speech, and
 201 by 30% for casual speech. A further motivation for introducing these duration
 202 classes is to contribute to future improvements of acoustic pronunciation
 203 modeling by estimating specific acoustic models for each class.

204 3.1. Method

205 In the alignment system used here, each acoustic phone model corre-
 206 sponds to a three state left-to-right hidden Markov model (HMM). Each
 207 phoneme is associated to one (or several) context-dependent acoustic model(s)
 208 corresponding to the three state left-to-right HMMs. The three states are
 209 assumed to model phone segment onset (first state), middle (second state)
 210 and end parts (third state). In forced alignment, each state is associated to
 211 at least one acoustic vector of 10 ms. As there are three states, the mini-
 212 mum duration of a phone segment amounts to 30 ms. Figure 2 illustrates
 213 forced alignment with a canonical (citation form) word pronunciation in a
 214 temporal reduction situation. From an ASR speech modelling perspective,
 215 temporal speech reduction may be captured by sequential pronunciation vari-
 216 ants with a smaller number of phonemes than the canonical pronunciation.
 217 Such productions, when aligned against canonical pronunciations, generally
 218 include one or several phone segments of a minimal duration (i.e. 30 ms
 219 here). These types of variants are considered to be the most problematic
 220 to ASR systems, since improper alignment results in poorer acoustic phone
 221 model accuracy.

222 Figure 3 gives an overview of how the ASR system can be used as an in-
 223 strument for linguistic purposes. Starting with manually transcribed speech
 224 at the word-level that is not time-aligned with the audio files, pronunciations
 225 that deviate from an expected norm may be located. The expected norm
 226 may be a citation form pronunciation or simply an a priori fixed minimum
 227 word or phone segment duration. The idea is to select portions of speech

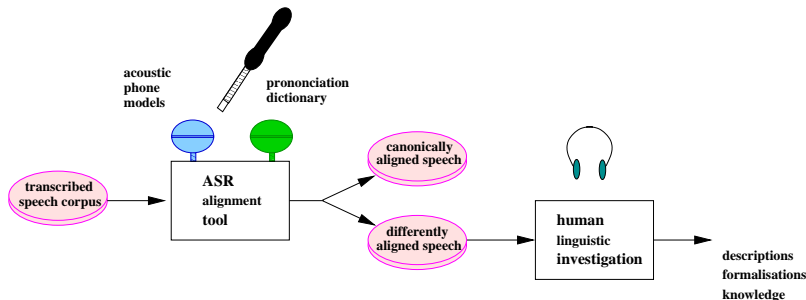


Figure 3: The automatic speech recognizer as an instrument to automatically select canonically and differently aligned subsets of speech deviating from expected representation. These subsets are of interest for more in-depth linguistic investigations.

that include a high proportion of such deviant forms (if they exist at all). Ideally, the adapted ASR tool might select the exact subset including all the deviant forms. In our case, the question of interest is temporal speech reduction, and therefore the number of deviant forms that are selected simply correlates with the number of low duration segments that is obtained by forced alignment. These subsets can then be studied more extensively by carrying out specific measurements. This approach may not only result in the improvement of the ASR model, but it may also enhance our knowledge of the linguistic phenomenon under investigation.

The proposed method aims at highlighting temporal reduction tendencies in large speech corpora via two different system configurations:

- In the first configuration, the alignment system makes use of a canonical pronunciation dictionary (see Table 4), which corresponds to the basic acoustic word model assumption of linear phoneme sequences as beads on a string. The approach consists of focusing on temporal structure mismatches

<i>word</i>	<i>canonical pronunciation</i>
le	lə
les	le lez(V)
maintenant	mɛ̃tənã

Table 4: Excerpt of the canonical pronunciation dictionary.

between the produced speech and the corresponding word models (i.e. HMM topologies of phonemic pronunciation models) as illustrated in Figure 2. The forced speech alignment technique matches the temporal structure of the

models onto the observed speech signal. Consequently, temporal reduction phenomena should correlate with higher rates of very short segments.

• A second configuration makes use of a pronunciation dictionary including optional schwas. Table 5 shows some lexical entries with a maximal number of pronunciations (including all possible phonemic segments as suggested either by the French writing conventions, or by pronunciation habits). The maximum length pronunciation thus includes all possible schwas, reflecting all mute-e occurrences within the word and potentially adding an epenthetic schwa at the end of a word-final closed syllable (without a graphemic presence of mute-e). Temporal reductions that result from schwa deletions can

<i>Word</i>	<i>Max. length pronunciation + variants</i>
Potential schwa: within word graphemic mute-e	
<u>le</u>	lə + l
<u>cela</u>	səla + sla
<u>devenu</u>	dəvəny + dəvny dvəny dvny
Potential schwa: closed-syllable word end	
<u>revanche</u>	ʁəvɑ̃ʃə + ʁəvɑ̃ʃ ʁvɑ̃ʃə ʁvɑ̃ʃ
<u>devenir</u> #	dəvənirə + dəvənir dəvnirə...

Table 5: Excerpt of the pronunciation dictionary for the schwa study. All possible variants are summed (+ sign) and arise from within word or from word-final optional schwas.

be measured via schwa deletion rates.

3.2. Speech corpora

For the purpose of our study, several important speech corpora were used, among which broadcast news (BN) and conversational telephone speech (CTS) corpora. The careful speech data-set stems from French broadcast news (BN) and corresponds to 360 hours of various radio and TV shows that were used for the *Technolange*-ESTER (Galliano et al., 2005) campaign. The casual speech data-set stems from the French telephone conversation (CTS) corpus and corresponds to 120 hours of LIMSI internal resources. French conversations from the corpora often took place between friends and/or family members, so the corpus therefore contains a highly casual speaking style. An additional French corpus (PFC) was used and provides different styles from the same set of users. The PFC (Phonologie du Français Contemporain, <http://www.projet-pfc.net/>) corpus is the result of an ambitious long-term project, initiated by French phonologists and phoneticians (cf. Durand,

271 Laks & Lyche, 2002, 2005). Several tens of varieties of the French language
272 from different regions across the French-speaking areas in the world have been
273 gathered. This amounts to the collection of data from hundreds of speakers
274 to enable large scale studies on linguistic and sociolinguistic variation found
275 in spoken French. For the present study, speech portions from ten regions
276 have been used, corresponding to ten hours of speech. The speech portions
277 were equally distributed between read speech and two sets of spontaneous
278 speech: supervised conversation and free conversation (the latter having a
slightly higher degree of casualness than the former one).

	French	
	# word tokens	duration
<i>Careful</i>	3600 k	360 h
<i>Casual</i>	1000 k	100 h
	English	
<i>Careful</i>	7200 k	720 h
<i>Casual</i>	25000 k	2300 h

Table 6: Corpus sizes for careful (BN) and casual (CTS) speech, for French (upper panel) and English (lower panel).

279
280 For comparison purposes, data from English corpora were added to the
281 French data. The English careful speech data-set included hundreds of hours
282 of broadcast news data for the DARPA *Rich Transcription 2004 Broadcast*
283 *News* evaluation (Nguyen et al., 2005), distributed by LDC (Linguistic Data
284 Consortium, <http://www.ldc.upenn.edu/Catalog/>). The casual speech data-
285 set stems from the Switchboard corpus (Godfrey, Holliman & McDaniel,
286 1992) and the more recent Fisher data (see LDC) including thousands of
287 hours of speech. In that corpus, telephone callers do not know each other and
288 are supposed to speak about assigned topics. Therefore, the speech, although
289 spontaneous, is less casual here than the speech in the French corpus. Each
290 corpus includes hundreds of male and female speakers.

291 3.3. Phone segment duration results

292 3.3.1. Casual versus careful speech styles

293 Figure 4 provides a line histogram of the segment proportions in the
294 French corpus as a function of segment duration (expressed in ms). Re-
295 sults are broken down for prepared and conversational speech styles. The

296 results show that the majority of segments are part of segment durations
 297 up until about 150 ms. Concerning prepared speech, the duration bin with
 298 the largest number of segments ($> 14\%$) corresponds to 60 ms. The casual
 299 speech distribution has by far the most segments ($> 18\%$) in the shortest du-
 300 ration bin of 30 ms, the total number of segments of 30 and 40 ms amounts
 301 up to more than 30% of the corpus. The distribution from the conversa-
 302 tional speech corpus shows a somewhat flattened distribution, which suggest
 303 that casual speech is characterized by more fast and more slow speech. A
 304 Kolmogorov-Smirnov test confirmed the significant difference between the
 305 casual and careful distributions ($D = 0.105$, $p < .0001$) More detailed analyses
 306 on a speaker-per-speaker basis confirmed that the trend generally holds for
 307 each individual speaker and that the observed flattening is therefore not a
 308 result of mere averaging pools of fast and slow speakers. The lengthened
 309 segments may be partly due to prosodically stressed items (cf. Adda-Decker
 310 et al., 2008), and partly due to hesitation phenomena.

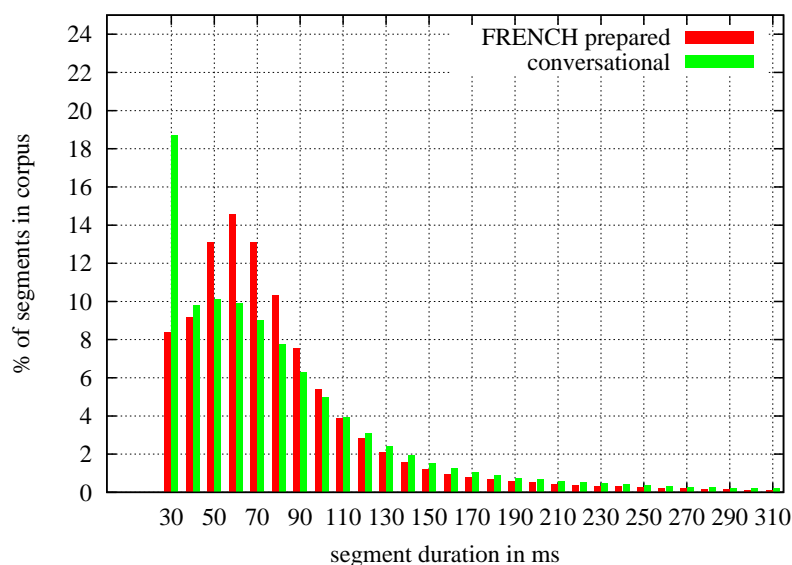


Figure 4: Line histogram of segment proportions from the corpus sets (cf. Table 6) as a function of segment durations for French. Journalistic prepared speech and conversational telephone speech styles are being compared.

3.3.2. French versus English

If casual speech generally produces a more flattened phone duration distribution, then it might be expected that this pattern should be observed irrespective of the corpus language. In other words, the observation should also hold for English. Figure 5 provides a line histogram of the segment proportions for English casual and careful speech as a function of segment duration.

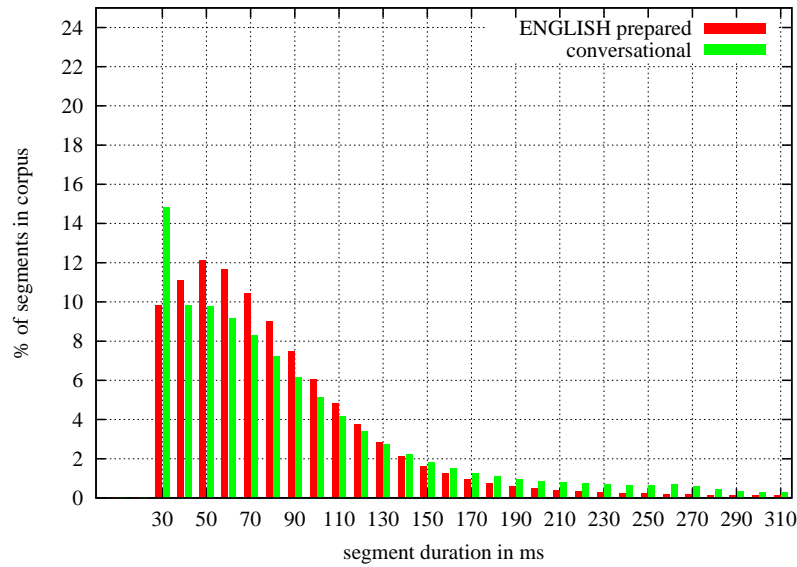


Figure 5: Line histogram of segment proportions from the corpus sets (cf. Table 6) as a function of segment durations for English. Journalistic prepared speech and conversational telephone speech are being compared.

Whereas the overall pattern slightly changes compared to the French results, the general tendency remains the same. There is a high proportion of segments in the minimum duration bin and a flattening of the conversational speech distribution, with more segments appearing in the longer segment durations (up until 310 ms). A Kolmogorov-Smirnov test confirmed the significant difference between the two speech styles ($D=0.084$, $p<.001$). Furthermore, it can be observed that the differences between careful and casual speech proportions up until 150 ms appear to be smaller in English than in French.

3.3.3. Comparing various casual speech styles

The comparisons between French and English were based on prepared journalistic speech and conversational speech styles taken from telephone conversations. To ascertain whether the tendencies found for spontaneous speech also hold for communicative settings other than telephone conversations, data from the French PFC corpus have been added. This corpus has the advantage of including a range of different speaking styles with read speech and face-to-face conversations, all produced by the same set of speakers.

Figure 6 shows the line histogram of segment proportions as a function of duration found for the PFC corpus (left-panel graph) and, for the reader's convenience, recalls the histogram for the previously shown French prepared and conversational corpus data (right-panel graph). The face-to-

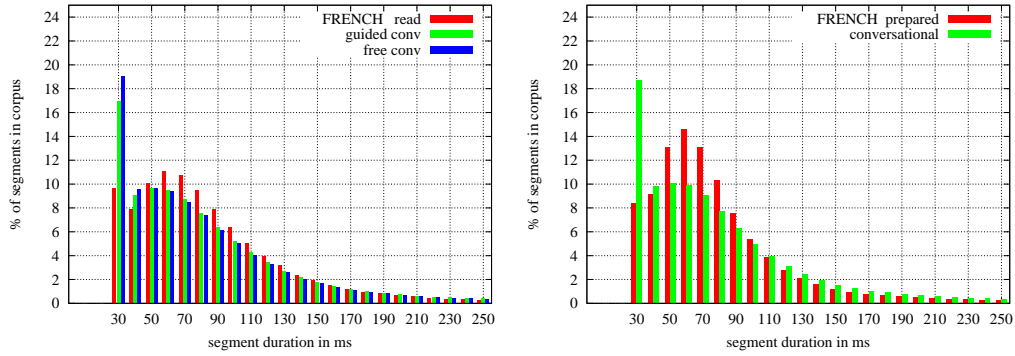


Figure 6: Segment duration proportions corresponding to read and conversational speech from the French PFC corpus (left) and journalistic prepared speech and conversational telephone speech (right).

face conversation distributions (guided and free) from the PFC corpus almost overlap with the conversational speech distribution, with quite similar casual speech effects. The read speech distribution was found to be significantly different from both the guided and free conversation distributions (Kolmogorov-Smirnov: $D = 0.089$, $p < .001$, and $D = 0.108$, $p < .001$ respectively). No significant difference was obtained between the PFC guided and free conversation distributions. Moreover, the read speech distribution exhibits a similar pattern as the prepared speech distribution on the right. However, it can be observed that there is a higher proportion of segments in the segment durations (up until 150 ms) for prepared journalistic speech as opposed to read PFC speech. This difference may very well be related

350 to time-pressure and speech rate performances for professional journalists,
 351 whereas it can be assumed that unprofessional readers speak more slowly.

352 3.3.4. Vowels versus consonants

353 The next question of interest is to know whether vowels or consonants
 354 exhibit similar duration patterns, and more specifically, which phonemes are
 355 the most prone to temporal reduction. As we have seen before, an obvi-
 356 ous candidate for high temporal reduction rates is the schwa (e.g. in high
 357 frequency function words such as *le*, *de*, *ne...*, ('the', 'of', 'not'...)). More-
 358 over, if it is true that function words are most prone to reduction, then the
 359 phonemes /l/ and /d/ should also be good candidates. Figure 7 shows overall
 360 proportions of consonants, vowels and sounds other than speech (this latter
 361 category includes real silences, but also breathing, hesitations and various
 362 other noises). The left-hand panel exhibits distributions taken from the pre-
 363 pared speech corpus, the right-hand one exhibits distributions taken from
 364 the spontaneous speech corpus. From the right-hand spontaneous speech
 365 distribution, it can be seen that vowels are only slightly more reduced and
 366 lengthened than consonants. This pattern holds for both careful and casual
 speaking styles.

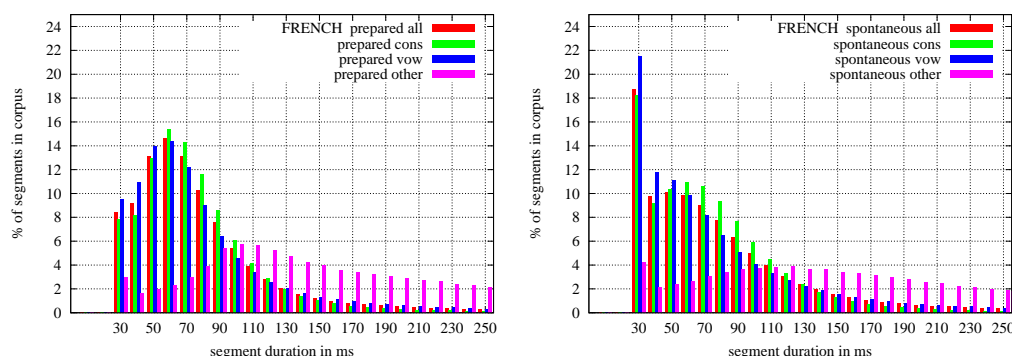


Figure 7: Segment duration distributions of vowels, consonants and other non-speech sounds (including silence) in prepared vs. spontaneous speech in French. For readability, each distribution sums up to 100%, even though the corpus is composed of 54% of consonants, 43% of vowels and 3% of other events (left) and of 50% of consonants, 43% of vowels and 7% of other events (right).

367

368 We will now be looking at duration phenomena for vowels and conso-
 369 nants separately. As was mentioned before, the segment proportions are
 370 represented as a function of the four duration classes, rather than continuous

371 segment durations. Turning to oral vowels, their duration distributions are
 372 shown in Figure 8. Given previous studies on the special status of the French
 373 schwa vowel mentioned earlier, it does not come as a surprise that the schwa
 374 vowel (the bar at the rightmost position) exhibits the highest figures in the
 375 short duration class, far above all the other vowels. Interestingly, the other
 376 French central vowel /ø/ (coded /eu/ in the figure, and near-homophone with
 377 a realized schwa) behaves in the opposite direction: it is actually one of the
 least reduced vowels. The duration patterns seen from the Figure, confirm

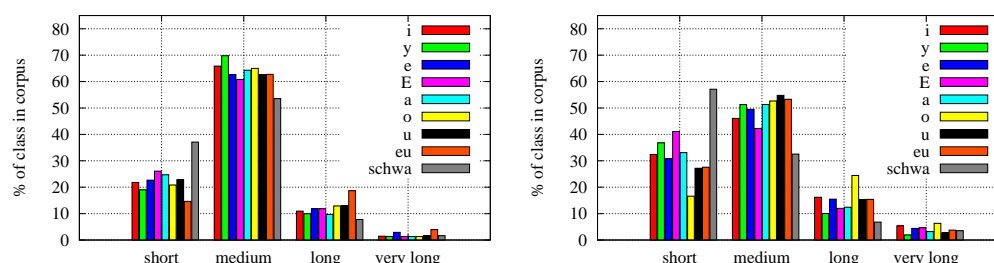


Figure 8: Proportions of short/medium/long/very long duration classes for French oral vowels. The left-hand panel shows the vowel distribution for careful speech distributions, the right-hand panel shows the vowel distribution for casual speech.

378 the special status of the French schwa vowel as an optionally (dis)appearing
 379 sound. Front and open vowels, such as /ε/ are good follow-up candidates in
 380 the temporal reduction hierarchy and tend to be temporally more reduced
 381 than back closed and rounded vowels.
 382

383 We now move on to duration patterns for consonants. Figures 9 and 10
 384 show duration distributions for the French voiceless and voiced fricatives
 385 respectively. Just like for the vowels, the results are shown for careful (left-
 386 hand panel) and casual speech (right-hand panel).

387 As can be seen from the Figures, the distribution patterns show that
 388 voiceless fricatives tend to be longer compared to the consonants' average
 389 duration (the bar at the rightmost position). In the *short* duration class the
 390 rates of the voiceless fricatives remain relatively low, particularly for /f/.

391 Although it is well-established that voiced consonants tend to have shorter
 392 durations than their voiceless counterparts, the /v/ exhibits a somewhat
 393 atypical behavior with respect to /z/ and /ʒ/, in that a very high rate of /v/

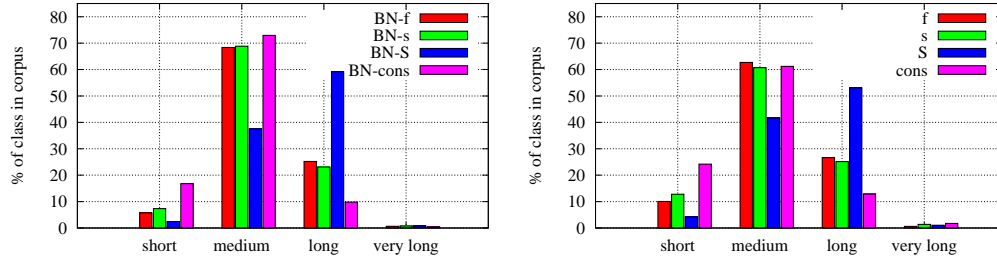


Figure 9: Proportions of short/medium/long/very long duration classes for French voiceless fricatives (/f/, /s/, /ʃ/) found in careful (left) and casual (right) speech styles.

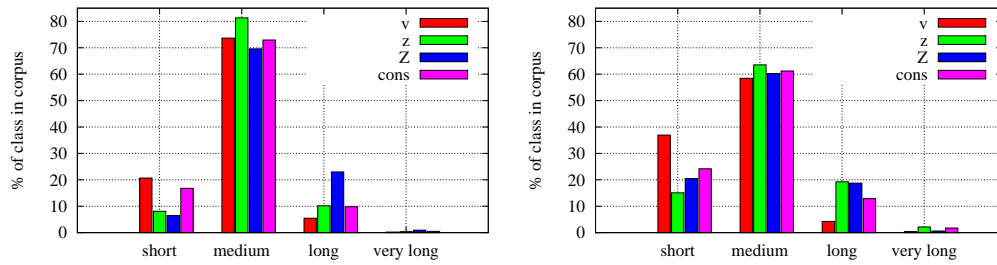


Figure 10: Proportions of short/medium/long/very long duration classes for French voiced fricatives (/v/, /z/, /ʒ/) found in careful (left) and casual (right) speech styles.

394 appear in the short duration class, with over 20% and 37% of segments for
 395 respectively careful and casual speech styles. In previous studies (see Adda-
 396 Decker et al., 2005), temporal reduction was already being observed for the
 397 function word *avec* /avek/ ('with'), that may be pronounced as something
 398 that is approximating [ɛg].

399 Finally, Figure 11 shows the results for the liquid and glide consonant
 400 classes in French. As opposed to the voiceless fricatives, liquids and glides
 401 tend to be rather short, given the increased proportions in the short du-
 402 ration class. The liquid /l/ which appears very often in frequent function
 403 words such as *le*, *la*, *les* ('the') appears to be most reduction-prone. For
 404 casual speech, the /ɥ/ glide also exhibits high reduction rates. These are

mainly stemming from the conversational specific words *suis* ('am') and *puis* ('then'). In particular the word sequence *je suis* ('I am') is frequently shortened to a pronunciation such as [ʃɥi] due to schwa-deletion and some complex consonantal assimilation processes.

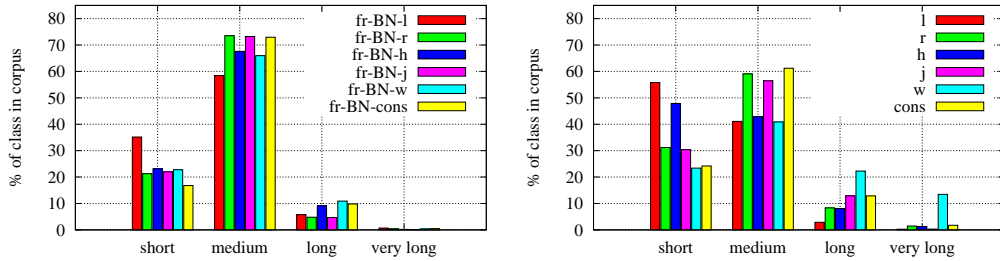


Figure 11: Proportions of short/medium/long/very long duration classes for French French liquids and glides found in careful (left) and casual (right) speech styles.

3.4. Optional schwa alignment

Instead of looking for evidence of temporal reductions in duration distributions from canonical alignments, the same objective can be pursued with a different tuning of the ASR model (see Adda-Decker, 2007, for more technical details). By adding sequential pronunciation variants to the pronunciation dictionary that includes optional segments (see Figure 5), the alignment system is free to keep or retrieve segments depending on their durations. It is clearly beyond the scope of the present article to provide an exhaustive overview of the implications of such adjustments to the model. However, we would like to give a flavour of the type of exploration one might conduct when focusing on French schwa realization and deletion.

Table 7 shows some schwa realization results concerning the French article *le* ('the') in the careful speech ESTER Corpus. Results highlight the contextual dependency of the schwa realization. Whereas the schwa is produced in almost 83% on average, there are left contexts where the rate approximates almost 100% (see, for instance, the rates for a left context of breathing (capital B in table 7) and in right unvoiced plosive /t/ context). On the other hand, a left vocalic context such as /u/ and the right voiced consonant /m/

context entails very low schwa production rates (14.3%). This specific context matches the idiom *tout le monde* [tuləmōdə] ('everyone'), which is most frequently produced as a bisyllabic word [tulmōd]. Globally, these results show that schwa production rates are about 90%, if the left context includes either breathing or a word ending with a closed syllable. A left vocalic context favours a schwa deletion. This is a good illustration of interesting temporal reduction phenomena such as can the case for the sequence *quai de Seine* ('dock of the Seine') that can be pronounced *quête saine* ('healthy quest').

cond.	%ə	%∅	nb.
#Cə	82.6	17.4	23950
B#Cə	92.7	7.3	2730
C#Cə	89.6	10.4	6480
V#Cə	76.7	23.3	12310
B#Cə#t	98.8	1.2	170
r#Cə#t	94.7	5.3	360
r#Cə#m	83.6	16.4	540
u#Cə#m	14.3	85.7	230

Table 7: Schwa realization (ə) and (complementary) deletion (∅) rates (expressed in %) of the French article *le* ('the') for the general pattern #Cə, and different left and right context conditions (#: word boundary, B: breathing, C: consonant, V: vowel). The total number of occurrences per condition (*nb*) are given in the last column.

4. Discussion

Whereas it may seem like a trivial matter to cite a number of examples of more or less severe reduction phenomena in day-to-day spoken language, the exact mechanisms that underly the processing of pronunciation variants that arise from these reductions are still relatively unknown. They are equally responsible for relatively high error rates in current ASR systems. More extensive descriptions are required to gain a better understanding of the type of pronunciation variation that both human listeners and automatic speech recognition systems are dealing with. By looking into transcription errors, it appears that a large number of ASR transcription errors in casual speech is due to a variety of "reduction phenomena" resulting in shorter productions. Given this observation, it seems straightforward to try and use

447 ASR systems as a tool to detect, quantify, and describe such phenomena
 448 using large spoken corpora. In the present paper, we focused on temporal
 449 patterns in French. Segmentation and labelling of speech regions that are
 450 prone to reduction phenomena were carried by using forced speech alignment
 451 in combination with canonical (full form) pronunciations. The forced speech
 452 alignment technique enabled us to localize and quantify phonemes that are
 453 particularly prone to temporal reduction. It must be borne in mind that
 454 upon employing forced alignment, an isolated occurrence of such a segment
 455 is not necessarily indicative of reduction. After all, an increase in duration
 456 of 50 ms will be relatively greater for a vowel that has a typical duration
 457 of 50 ms than for one that has a typical duration of 150 ms (see Campbell,
 458 1992). Nevertheless, the larger the number of contiguous minimum duration
 459 segments, the stronger the hypothesis of an actual temporal reduction.

460 In a series of explorations in French, phone segment duration distribu-
 461 tions were computed for a number of different speaking styles and speaker
 462 populations, ranging from careful speech (i.e., read speech, prepared journal-
 463 istic speech) to casual speech (i.e., telephone conversations and face-to-face
 464 interviews). With regard to more refined comparisons between French vow-
 465 els and consonants, the segment duration distributions were examined using
 466 four duration classes.

467 Our results showed that casual speech has a flatter duration distribution
 468 than careful speech, with an increase of more than 10% in the proportion of
 469 short segments, but also a larger number of longer segments. This general
 470 trend has been confirmed for English, even though the increase of short seg-
 471 ments was limited to about 6%. This may be partly due to the fact that
 472 the corpora used for English included less casual telephone conversations.
 473 Future research should aim to investigate these measured differences with
 474 respect to known prosodic syllable/stress-timing differences between French
 475 and English (see, e.g., Cutler, Mehler, Norris & Segui, 1986). Moreover, ob-
 476 servations also remain stable by shifting from prepared and telephone speech
 477 to the PFC corpus, where the same speaker population produced both read
 478 and spontaneous face-to-face speech. Next, duration patterns in vowels and
 479 consonants were examined. Vowels are slightly more shortened than conso-
 480 nants and, not surprisingly, the schwa vowel, that is notoriously known for
 481 being deleted in French, exhibits the highest figures of short segments, far
 482 above all other vowels. As for reduction patterns in consonants, it was shown
 483 that for voiced fricatives, the /v/ exhibits an unexpected duration pattern,
 484 in that it is relatively short compared to the other voiced fricatives. The

485 liquid /l/ which appears very often in frequent function words is most prone
486 to shortening. In casual speech, the /ɥ/ glide also exhibits high reduction
487 rates due to conversation-specific words *suis* ('am'), *puis* ('then'). Finally, a
488 schwa-specific exploration using forced alignment and an adapted pronunci-
489 ation dictionary showed regularities of schwa production before pauses and
490 breathing, and schwa deletions in idiomatic expressions such as *tout le monde*
491 ('everyone').

492 By using forced alignment to quantify temporal reduction phenomena in
493 French, we hope to have demonstrated how ASR systems may serve as a tool
494 to systematically investigate variations that occur in the speech signal across
495 different speaking styles. Hopefully, the present results will shed some new
496 light on the intrinsically complex nature of temporal processes in speech. In
497 future work, we plan to refine the present approach and to further extend
498 the analysis of the alignment results. Studying linguistic phenomena from
499 an ASR perspective using large corpora might also give us some clues about
500 the encoding of information in speech. The speech signal is endowed with
501 fine phonetic detail and features that the human listener seems to rely on
502 in the face of ambiguity and noise. The perspectives available through an
503 ASR approach are manifold. For researches working in the domain of ASR,
504 the ultimate goal is to uncover the generic rules to generate pronunciation
505 variants, even for rarely observed or unobserved words, for which variants
506 cannot be estimated statistically. The framework developed should also help
507 to describe and quantify more or less well known linguistic phenomena on
508 phonemic and lexical levels, which is of relevance to linguists and cognitive
509 scientists alike.

510 5. Acknowledgement

511 Parts of the research reported in this article have been funded by grants
512 from the CNRS, ANR, Digiteo, and Quaero, awarded to the first author, and a
513 grant from the Luxembourgish Fondation Nationale de la Recherche, awarded
514 to the second author. We would like to thank Lori Lamel, Jean-Luc Gauvain,
515 and Gilles Adda for their help and fruitful discussions during the preparation
516 of this paper. We are greatly indebted to three anonymous reviewers for their
517 constructive criticisms and suggestions on an earlier version of the paper.

518 References

- 519 Adda-Decker, M. [http://ling.upenn.edu/phonetics/workshop.](http://ling.upenn.edu/phonetics/workshop/), Gendrot C. &
520 Nguyen N., (2008). Contributions du traitement automatique de la pa-
521 role à l'étude des voyelles orales du français. *Traitement Automatique des*
522 *Langues*, 49(3), 13-46.
- 523 Adda-Decker, M. (2007). Problèmes posés par le schwa en reconnaissance
524 et en alignement automatiques de la parole. In: *Actes des 5es Journées*
525 *d'Études Linguistiques de Nantes*, Nantes, 211-216.
- 526 Adda-Decker, M., Boula de Mareüil, P., Adda, G., Lamel, L., (2005). In-
527 vestigating syllabic structures and their variation in spontaneous French.
528 *Speech Communication*, 46, 119-139.
- 529 Adda-Decker, M. & Lamel, L. (2005). Do speech recognizers prefer female
530 speakers? In: *Proceedings of InterSpeech*, Lisbon.
- 531 Adda-Decker, M. & Lamel, L., (1999). Pronunciation variants across system
532 configuration, language and speaking style. *Speech Communication*, 29, 83-
533 98.
- 534 Campbell, N. (1992). Segmental elasticity and timing in Japanese speech. In:
535 Tohkura, Vatikiotis-Bateson, and Sagisaka, Speech perception, production
536 and Linguistic Structure. IOS Press, Amsterdam, Washington, Oxford,
537 pp.403-418.
- 538 Cole, R., Oshika, B., T., Noel, M., Lander, T., Fanty, M. (1994). Labeler
539 Agreement in Phonetic Labeling of Continuous Speech. In: *Proceedings. of*
540 *ICSLP*, 2, 2131-2134.
- 541 Cutler, A., Mehler, J., Norris, D., & Segui, J., 1986. The Syllable's Differing
542 Role in the Segmentation of French and English. *Journal of Memory and*
543 *Language*, 25, 385-400.
- 544 Dausès, A., (1973). Études sur l'e instable dans le français familier. Niemeyer
545 Verlag. Tübingen.
- 546 Dilley, L., Pitt, M. (2007). A study of regressive place assimilation in spon-
547 taneous speech and its implications for spoken word recognition. *Journal*
548 *of the Acoustical Society of America*, 122, pp.2340-2353.

- 549 Duez, D., (2003). Modelling Aspects of Reduction and Assimilation in Spon-
 550 taneous French Speech. In: *Proceedings of the IEEE-ISCA Workshop on*
 551 *Sponrtaeous Speech Processing and Recognition*, 2003. Tokyo.
- 552 Durand, J., Laks B. & Lyche C. (2002). La phonologie du français contempo-
 553 rain: usages, variétés et structure. In: C. Pusch & W. Raible (eds.) Roman-
 554 istische Korpuslinguistik - Korpora und gesprochene Sprache / Romance
 555 Corpus Linguistics - Corpora and Spoken Language. Tübingen: Gunter
 556 Narr Verlag, pp. 93-106.
- 557 Durand, Jacques, Laks B. & Lyche C. (2005). Un corpus numérisé pour la
 558 phonologie du français. In G. Williams (ed.) *La linguistique de corpus*.
 559 Rennes: Presses Universitaires de Rennes. pp. 205-217.
- 560 Elman, J. & McClelland, J.M. (1988). Cognitive penetration of the mecha-
 561 nisms of perception: Compensation for coarticulation of lexically restored
 562 phonemes. *Journal of Memory and Language*, 27, 143-165.
- 563 Ernestus M. (2000). Voice assimilation and segment reduction in casual
 564 Dutch, a corpus-based study of the phonology-phonetics interface. Utrecht:
 565 LOT.
- 566 Fougeron, C., Goldman, J.-P. & Frauenfelder, U.H., (2001). Liaison and
 567 schwa deletion in French: an effect of lexical frequency and competition.
 568 In: *Proceedings of Eurospeech* , Aalborg, 639-642.
- 569 Galliano et al. (2005). The Ester Phase II Evaluation Campaign for the Rich
 570 Transcription of French Broadcast News. In: *Proceedings of InterSpeech*.
- 571 Ganong, W.F. (1980). Phonetic categorization in auditory word perception.
 572 *Journal of Experimental Psychology: Human Perception & Performance*,
 573 6, 110-125.
- 574 Gauvain J.-L., Lamel L.F., Adda G. & Adda-Decker M. (1994). Speaker-
 575 Independent Continuous Speech Dictation. *Speech Communication*, 15(*).
- 576 Gauvain, J.-L., Adda, G. , Adda-Decker, M., Allauzen, A., Gendner, V.,
 577 Lamel, L., Schwenk, H. (2005). Where Are We in Transcribing French
 578 Broadcast News? In *Proceedings of InterSpeech*, Lisbon.

- 579 Gendrot C. & Adda-Decker M., (2005). Impact of duration on F1/F2 for-
580 mant values of oral vowels: an automatic analysis of large broadcast news
581 corpora in French and German. In: *Proc. of InterSpeech*, Lisbon.
- 582 Godfrey J., Holliman E. & McDaniel J. (1992). Switchboard: telephone
583 speech corpus for research and Development. In: *Proceedings of IEEE-*
584 *Icassp*.
- 585 Greenberg, S. & Chang, S., 2000. Linguistic dissection of Switchboard-Corpus
586 Automatic Speech Recognition Systems. In: *Proc. ISCA-ITRW Workshop*
587 *on ASR*, Paris, pp. 195-202.
- 588 Greenberg, S., Carvey, H., Hitchcock, L., Chang, S. (2003). Temporal prop-
589 erties of spontaneous speech – a syllable-centric perspective. *Journal of*
590 *Phonetics*, 31, 465-485.
- 591 Hosom, J.-P. (2009). Speaker-independent phoneme alignment using
592 transition-dependent states. *Speech Communication*, 51(4), 352-368.
- 593 Jurafsky D., Bell A., Gregory M. & Raymond W.D., (2001). Probabilistic
594 relations between words: Evidence from reduction in lexical production
595 in *Frequency and the Emergence of Linguistic Structure*, Bybee & Hopper
596 eds. pp. 229-254, John Benjamins.
- 597 Lefèvre F., Gauvain J.-L., Lamel L.F., (2005). Genericity and portability
598 for task-dependent speech recognition. *Computer Speech and Language*,
599 19:345-363, 2005.
- 600 Nakamura M., Furui S. & Iwano K. (2006). Acoustic and Linguistic Char-
601 acterization of Spontaneous Speech Masanobu, *ISCA workshop on Speech*
602 *Recognition and Intrinsic Variation*, Toulouse, France.
- 603 Nguyen L., Abdou S., Afify M., Makhoul J., Matsoukas S., Schwartz R., Xi-
604 ang B. Lamel L., Gauvain J.L., Adda G., Schwenk H., Lefèvre F., (2005).
605 The 2004 BBN/LIMSI 10XRT English Broadcast News Transcription Sys-
606 tem. In *IEEE-Icassp 2005*.
- 607 Prasad R., Matsoukas S., Kao C.L., Ma J., Xu D.X., Gauvain J.-L., Lamel
608 L., Schwenk H., Adda G. & Lefèvre F., (2005). The 2004 BBN/LIMSI
609 20xRT English Conversational Telephone Speech Recognition System. In:
610 *Proceedings of. InterSpeech*, Lisbon.

- 611 Samuel, A.G. & Pitt, M.A. (2003). Lexical activation (and other factors) can
612 mediate compensation for coarticulation. *Journal of Memory and Lan-*
613 *guage*, 48, 416-434.
- 614 Schuppler B., Ernestus M., Scharenborg O. & Boves L. (2008). Corpus of
615 Dutch Spontaneous Dialogues for Automatic Phonetic Analysis. In: *Pro-*
616 *ceedings of Interspeech*, Brisbane, pp. 1638-1641.
- 617 Snoeren, N., Hallé, P. & Segui, J. (2006). A voice for the voiceless : Produc-
618 tion and perception of assimilated stops in French. *Journal of Phonetics*,
619 34, 241-268.
- 620 Shriberg, E. (1994). Preliminaries to a Theory of Speech Disfluencies. PhD
621 thesis, University of California, Berkeley.
- 622 Strik H., Binnenpoorte D. & Cucchiariini C. (2005). Multiword Expressions
623 in Spontaneous Speech: Do we really speak like that? In: *Proc. of Inter-*
624 *Speech*, Lisbon, pp.1161-1164.
- 625 Strik H., Elffers A., Bavcar D. & Cucchiariini C., (2006). Half a Word is
626 Enough for Listeners, but Problematic for ASR. In: *Proceedings of ISCA*
627 *workshop on Speech Recognition and Intrinsic Variation*, Toulouse, France.
- 628 Strik, H., Cucchiariini, C., (1999). Modelling pronunciation variation for ASR:
629 A survey of the litterature. *Speech Communication*, 29,, 225-246.
- 630 Van Son, R.J.J.H., Pols & L.C.W., (2003). An Acoustic Model of Commu-
631 nicative Efficiency in Consonants and Vowels taking into Account Con-
632 text Distinctiveness. In: *Proceedings of the 15th ICPhS*, Barcelona 2003,
633 pp. 2141-2143.
- 634 Tseng S.-C., Features of Contracted Syllables of Spontaneous Mandarin, In:
635 *Proceedings of InterSpeech*, Lisbon.
- 636 Verney Pleasants, J., (1956). Études sur l'e muet, timbre, durée, intensité,
637 hauteur musicale. Klincksieck, Paris.